

NUMERICAL METHODS

—
BOOTH

THIRD EDITION



NUMERICAL METHODS — BOOTH

TH
321

£2

NUMERICAL
METHODS

NUMERICAL METHODS

NUMERICAL METHODS

By

Donald

ANDREW D. BOOTH, D.Sc.

*Dean of the College of Engineering,
University of Saskatchewan,
Interdisciplinary Professor of Autometrics,
Western Reserve University, Cleveland*

THIRD EDITION

LONDON
BUTTERWORTHS
1966

ENGLAND: BUTTERWORTH & CO. (PUBLISHERS) LTD.
LONDON: 88 Kingsway, W.C.2

AUSTRALIA: BUTTERWORTH & CO. (AUSTRALIA) LTD.
SYDNEY: 20 Loftus Street
MELBOURNE: 473 Bourke Street
BRISBANE: 240 Queen Street

CANADA: BUTTERWORTH & CO. (CANADA) LTD.
TORONTO: 1367 Danforth Avenue, 6

NEW ZEALAND: BUTTERWORTH & CO. (NEW ZEALAND) LTD.
WELLINGTON: 49/51 Ballance Street
AUCKLAND: 35 High Street

SOUTH AFRICA: BUTTERWORTH & CO. (SOUTH AFRICA) LTD.
DURBAN: 33/35 Beach Grove

U.S.A.: BUTTERWORTH INC.
WASHINGTON, D.C.: 20014: 7300 Pearl Street

First Edition	1955
Second Edition	1957
Second Impression	1958
Third Edition	1966

511-4 B00

CLASS	511.4
VOL.	3
COPY	
SUPPLIER	Philip
REC'D	22.4.71
ACCESS	161

©
Butterworth & Co. (Publishers) Ltd.
1966

Set in Monotype Baskerville type

Made and printed by offset in Great Britain by
William Clowes and Sons, Limited, London and Beccles

PREFACE TO THIRD EDITION

REFLECTING the advances in the subject since the first edition, much new material has been added to the text. This new material takes account, not only of experience gained in using the earlier edition as a class text book, but also of our observations of the virtues of different methods on modern computing machines.

A. D. B.

Saskatoon

PREFACE TO FIRST EDITION

THE present book has grown out of the series of lectures given by the author to final honours B.Sc. mathematics students at Birkbeck College, London. The purpose of the course has been, not so much to instruct in the detailed tedium of actual calculation, but rather to give an understanding of the basic principles upon which such analyses rest.

Some existing books on numerical analysis lay much stress upon the detailed form in which a given procedure is to be laid out; it has been my experience that such concentration upon actual numbers obscures the underlying mathematical basis upon which the work rests, and so such tabulations are almost entirely absent from this book. Where they are given, as in Chapter 7, they illustrate the sort of behaviour which will be encountered in a calculation rather than any detailed form of layout.

Were these didactic points the sole reason there would be little justification for a new book on computation; a far more important consideration lies in the growth, during the past decade, of the science and art of programming for an automatic digital calculator. The classical methods of hand calculation are, to a greater or less extent, unsuitable for the modern machines, and only by having a thorough knowledge of the underlying mathematical principles, is the programmer likely to make effective use of the new tools.

At Birkbeck College Computational Laboratory the teaching of numerical methods has been accompanied by the actual use of an

PREFACE

automatic calculator, and demonstrations of such things as differencing and the solution of differential equations have been carried out by the machine and not by the student. Perhaps not unnaturally, this has proved more popular than the old method.

A book of this kind must always owe much to the work of previous authors; it is pleasant to acknowledge the help which the author derived from Freeman's 'Actuarial Mathematics' and from the classical 'Calculus of Observations' of Whittaker and Robinson, both of which were practically the only available works during the 1930's. In more recent times the paper on 'Difference and Associated Operators' by W. G. Bickley may be mentioned as having particular influence.

Finally it gives me particular pleasure to acknowledge the help of my wife both in making the book more readable than might otherwise have been the case, and also, in collaboration with J. P. Cleave, B.Sc., for checking the examples given in Chapter 7.

A. D. B.

Fenny Compton

AUTHOR'S NOTE

Advantage has been taken of the need for a second edition to correct small errors which were present in the earlier version, and the author wishes to express his thanks to reviewers and others who drew his attention to these.

A. D. B.

Fenny Compton

CONTENTS

	PAGE
PREFACE	v
1. THE NATURE AND PURPOSE OF NUMERICAL ANALYSIS .	1
2. TABULATIONS AND DIFFERENCES	7
3. INTERPOLATION	12
4. NUMERICAL DIFFERENTIATION AND INTEGRATION . .	30
5. THE SUMMATION OF SERIES	54
6. THE SOLUTION OF ORDINARY DIFFERENTIAL EQUATIONS .	59
7. SIMULTANEOUS LINEAR EQUATIONS	76
8. PARTIAL DIFFERENTIAL EQUATIONS	124
9. NON-LINEAR ALGEBRAIC EQUATIONS	156
10. APPROXIMATING FUNCTIONS	178
11. FOURIER SYNTHESIS AND ANALYSIS	187
12. INTEGRAL EQUATIONS	197
SELECT BIBLIOGRAPHY	205
NAME INDEX	209
SUBJECT INDEX	211

THE NATURE AND PURPOSE OF NUMERICAL ANALYSIS

1.1 HISTORY

ALTHOUGH numerical analysis is considered by some to be a subject of recent origin and development this is not, in fact, so. Dealing, as it does, with the derivation of results in the form of *numbers*, the numerical analyst is really the lineal descendant of the first caveman who enumerated the number of his wives by putting them into one : one correspondence with the fingers of his hand.

Even in its more modern aspects the subject is antique; thus a primary activity of the scientists of Babylon was the construction of mathematical tables. An example is extant, dating from about 2000 B.C., which contains on a tablet the squares of the numbers 1-60. Another tablet records the eclipses going back to 747 B.C., so that astronomical calculation formed a part of the activity of these early numerical analysts.

The ancient Egyptians, too, were energetic numerical analysts. They constructed tables whereby complex fractions could be decomposed into the sum of simpler forms with unit numerators, and invented the method of *false position* (see Chapter 9, section 9.3) for the solution of non-linear algebraic equations.

Passing to the Greek mathematicians we find Archimedes, in about 220 B.C., approximating the value of π and describing it as less than $3\frac{1}{7}$ but greater than $3\frac{1}{9}$. Heron the elder, in about 100 B.C., made use

of the iterative process: $\sqrt{a} \sim \frac{1}{2} \left(x_n + \frac{a}{x_n} \right)$ which is usually ascribed to Newton, and the Pythagorean school considered the summation of the series $(1 + 2 + 3 \dots)$. Diophantus, about A.D. 250, apart from his better known work on indeterminate equations, was responsible for a process for the arithmetical solution of quadratic equations.

The Hindus were the creators of our modern arithmetic notation—usually called Arabic—and devised the method of checking the correctness of an arithmetic calculation known as ‘casting out nines’.

Mohammed ibn Musa Al-Khowarizmi was the first Arab arithmetician and was responsible, around A.D. 820, for the systematization of computational processes. He gave the value $\pi = 62832/20000$ and was active in the preparation of astronomical tables. Abul Wefa

(A.D. 960) devised a method for the computation of tables of sines and gave the value of $\sin(\frac{1}{2}^\circ)$ correct to nine decimal places; he also used the *tangent* and calculated a table of this function.

Jumping to the seventeenth century, it is interesting to note that Napier's first table of logarithms was produced before the use of exponents was current, and that his 'logarithm' differs from any in current use since:

$$\text{Napierian log } x = 10^7 \log_e (10^7/x).$$

In 1614 Napier published his *Mirifici logarithmorum canonis descriptio* and, posthumously, his *Mirifici logarithmorum canonis constructio* in 1619. Briggs, only slightly later in 1624, produced his *Arithmetica logarithmica*, which contains the logarithms, to 14 places of the numbers 1–20,000 and 90,000–100,000. Vlacq produced, in 1628, a table which is still fundamental of the 14-place logarithms of the numbers 1–100,000. The first authoritative publication of the logarithms of trigonometric functions was made at about the same time (1620), by Gunter, who invented the words 'cosine' and 'cotangent', and was responsible for the so-called 'Gunter's chain'.

In the nineteenth century there occurred one of the triumphs of numerical analysis, the simultaneous prediction by Adams and le Verrier in 1845, of the existence and position of the planet Neptune. This century saw also the rise and development of automatic calculating machinery, from the crude desk multiplier of Thomas de Colmar to the almost unmodified Brunsviga of the present day, the Hollerith punched card census calculator, and the difference and analytical engines of Charles Babbage.

Not until the end of the 1930's did the fully automatic calculators begin to come into use, and since the late 1940's there has been a revolution and renaissance in numerical analysis. New methods have been developed and problems which could not previously have been contemplated, even for a life's work, are now solved in hours. It is perhaps dangerous to quote examples, but outstanding achievements are the calculation of π and e to 100,000 decimals⁽¹⁾ which took an I.B.M. 7090 computer just under 8½ h. The demonstration of the primeness of the Mersenne number $2^{1279} - 1$ on S.W.A.C. in 13 min. 25 sec., may also be cited as a noteworthy achievement.

1.2 THE TOOLS OF ANALYSIS (HAND)

From the classical standpoint of the individual numerical analyst the tools of computation are:

- | | |
|-------------------------------|------------------------------|
| (1) Tables of formulae | (3) A desk calculator |
| (2) Tables of function values | (4) Pencil, paper and rubber |

Few investigations are of such a fundamental nature that they make no use of existing mathematical knowledge; probably the most common table of formulae in use is a list of integrals. Four standard works may be mentioned:

- (1) PEIRCE, B.O., 'A short table of integrals,' Ginn, Boston (1929)
- (2) DWIGHT, H. B., 'Tables of integrals and other mathematical data,' Macmillan, New York (1934)
- (3) DE HAAN, D. B., 'Nouvelles tables d'intégrales définies, reprinted Stechert, New York (1939)
- (4) 'Interpolation and allied tables,' H.M. Stationery Office (1956)

The first two volumes are chiefly concerned with indefinite integrals, and the third exclusively with definite integrals. The last booklet contains most of the useful formulae for interpolation.

Tables of function values are almost too numerous to mention. For 4- or 5-figure accuracy there are the classical:

JAHNKE-EMDE, 'Tafeln höherer Funktionen,' Teubner (4th edn.), Leipzig (1948)

EMDE, 'Tafeln elementarer Funktionen,' Teubner, Leipzig (1940) which, besides giving numerical values, contain useful graphs and formulae. Other moderate accuracy collections are:

'Mathematical Tables from the Handbook of Chemistry and Physics,' Chemical Rubber Publishing Co. Cleveland (1946)

DALE, J. B., 'Five-figure Tables of Mathematical Functions,' Arnold, London (1937)

DWIGHT, H. B., 'Mathematical Tables,' McGraw-Hill, New York (1941)

More accurate tables (6 or more decimal digits) are:

'Chambers's Seven-figure Mathematical Tables,' Chambers, London (1937)

'Chambers's Six-figure Mathematical Tables' (2 vols.) (Ed. Comrie) London (1948–9)

'Barlow's Tables of Squares, Cubes and Reciprocals' (Ed. Comrie) Spon (1941)

For more specialist tables reference can be made to the monumental:

FLETCHER, A., MILLER, J. C. P., ROSENHEAD, L. and COMRIE, L. J., 'Index of Mathematical Tables,' 2nd edn., Addison-Wesley, Mass. (1962)

It may be supplemented by reference to 'Mathematical Tables and other Aids to Computation' (M.T.A.C.) which maintains a cumulative description of new tables.

Numerous desk calculators are now available but none can be said to possess such excellence as to satisfy all felt wants and to outshine the others. Our experience recommends the Brunsviga and Madas machines in the hand-operated range, and the Marchant, Madas and Mercedes-Euclid amongst those electrically operated. We shall not attempt any description of the means of using these machines, since a short time with an experienced operator and the machine will do more than many pages of words in this direction.

Regarding pencil and paper, we may remark that foolscap paper ruled with faint $\frac{1}{4}$ in. squares seems convenient for most general purposes. Pencils should be soft, B is suitable, and the rubber should not be one degraded by age and use!

1.3 THE TOOLS OF ANALYSIS (AUTOMATIC)

It is hoped that amongst the readers of this book there will be many who have access to one of the automatic digital calculators which become daily more easily available.

Fortunately most of the available automatic machines are provided with an autocoding system such as Fortran or Algol. The development of these autocodes has revolutionized the preparation of arithmetic problems for computation. The compiling programmes can, to a large extent, detect careless errors in preparation and the very extensive facilities for floating point arithmetic, which are built into modern autocodes, make the careful experimental arithmetic of the immediate past almost unnecessary.

1.4 PRECISION, ACCURACY AND ERRORS

In planning a calculation the three factors detailed in the heading must always be considered. First, in any calculation using data obtained by physical measurement, the inherent precision of the data themselves must be examined. If no figures for experimental errors are presented and these are not easily obtainable, a knowledge of the experimental technique may give a clue. Measurements of length are rarely accurate to better than 1/10 per cent, measurements of weight often attain 1/10,000 per cent. Electrical measurements are frequently of precision as low as 5 per cent. These circumstances should be taken into account at the planning stage, and a rough working rule is to calculate to two places more than those given by the data.

The accuracy of the calculation (excluding errors of the careless type) will depend on the numerical process involved. Additions and subtractions neither increase nor decrease the precision of the data;

multiplications and divisions, however, lead to round off procedures and thus to an overall decrease in accuracy. Hand calculations are seldom of such length as to cause trouble from the growth of round off errors, but with automatic calculators the situation is different. Thus a typical matrix inversion may lead to over 10,000 multiplications, and this, in turn, to a rounding error which has a probable value of the order of 100 units in the last place.

Errors are of two main types, mathematical and human. The former may result from the use of approximations, which are inevitable consequences of the use of discrete processes to represent continuous ones. The latter class of error should be avoidable, at least in the long run, by the provision of adequate checks.

In manual calculation it has long been a platitude that a person should never check his own work and that, if possible, the same method should not be used. This has tended to become forgotten in connection with the use of automatic digital calculators, and it is frequently suggested that because the machine has produced the same answer twice in succession it must be correct. Unfortunately most machines suffer at times from 'pattern sensitivity', that is they will work with complete accuracy on all numbers except *one*. Under these circumstances the same *wrong* answer can be produced as often as required, and the only valid check is a completely different computing routine.

In practice it is often possible to check the results of a long series of calculations by some completely external means, such as differencing (see Chapter 3, section 3.2), or, alternatively, from a knowledge of the observations with which the calculations are intended to agree. When this is not so (with an automatic digital calculator) a good plan is to repeat the calculation after an interval of several days, since few, if any, of these machines survive such a period without servicing and this almost always varies any pattern sensitivity which may be present.

Formulae may be used in different ways and with differing resultant accuracy. Thus, we may calculate

$$\sin \theta \approx \theta - \theta^3/6 + \theta^5/120 \quad \dots (1.4.1)$$

in two ways. In the first the terms are formed separately and then added. In this event the round off *may* be equal to 1.5 in the last place. On the other hand, by forming:

$$[(\frac{1}{120}\theta^2 - \frac{1}{6})\theta^2 + 1]\theta \quad \dots (1.4.2)$$

the greatest error will be 0.5 in the last place. This example is also instructive in that it illustrates an efficient means of calculating a

polynomial. Thus a direct calculation of equation 1.4.1 involves two divisions, four multiplications and two additions or subtractions, whereas in equation 1.4.2 the multiplications are reduced to three.

Considerations such as those mentioned above should always precede any numerical calculation and might be multiplied indefinitely. Some pointers will be given at appropriate places in the text, but pencil and paper analysis can only be learned by long practice and, in our experience, no two expert analysts agree upon the best detailed layout for any particular case. For this reason the details will be left to the reader and such numerical examples as appear are illustrations of such things as convergence, rather than of computational layout.

REFERENCE

- (1) SHANKS, D. and WRENCH, J. W., *Maths Comput.*, 16 (1962) 76

TABULATIONS AND DIFFERENCES

2.1 THE NATURE OF TABULATED FUNCTIONS

THE differential calculus had its origin in a consideration of the mode of variation of a function $y = f(x)$ with the augment x . In the process of defining the differential coefficient, $\frac{dy}{dx}$, it is necessary to consider the limit of a finite difference ratio:

$$\frac{f(x + \delta x) - f(x)}{\delta x}$$

as $\delta x \rightarrow 0$. When a function is represented by a set of numerical values contained in a table, it is natural to consider the analogues of the differentials dy and dx which can be deduced immediately from the tabular values. Suppose that a function u_x is defined by a table:

x	u_x
0	u_0
1	u_1
2	u_2
\vdots	
n	u_n

Then, corresponding to dx , we have $(1 - 0)$, $(2 - 1)$, $(n - \{n - 1\})$ all of which are equal to unity. And corresponding to dy we have the differences $(u_1 - u_0)$, $(u_2 - u_1)$, . . . $(u_n - u_{n-1})$.

Since the interval of tabulation (*i.e.* 1 in this case) is by no means always unity, it is customary to represent this by the symbol $\Delta(x)$, and in a like manner: $u_{m+1} - u_m$ is represented by $\Delta(u_m)$.

Now just as it is possible to proceed to differential coefficients higher than the first, so can the differences of a tabulated function be extended thus:

x	$f(x)$	Δf	$\Delta^2 f$	$\Delta^3 f$
0	u_0	$\Delta u_0 = u_1 - u_0$	$\Delta^2 u_0 = \Delta u_1 - \Delta u_0 = u_2 - 2u_1 + u_0$	$\Delta^3 u_0 = \Delta^2 u_1 - \Delta^2 u_0 = u_3 - 3u_2 + 3u_1 - u_0$
1	u_1	$\Delta u_1 = u_2 - u_1$	$\Delta^2 u_1 = \Delta u_2 - \Delta u_1 = u_3 - 2u_2 + u_1$	$\Delta^3 u_1 = \Delta^2 u_2 - \Delta^2 u_1 = u_4 - 3u_3 + 3u_2 - u_1$
2	u_2	$\Delta u_2 = u_3 - u_2$	<i>etc.</i>	<i>etc.</i>
3	u_3	$\Delta u_3 = u_4 - u_3$		
4	u_4	$\Delta u_4 = u_5 - u_4$		
5	u_5	$\Delta u_5 = u_6 - u_5$		

The differences on the first line, namely, Δu_0 , $\Delta^2 u_0$, $\Delta^3 u_0$ *etc.* are referred to as leading differences.

2.2 SOME ACTUAL TABLES

Before proceeding to a consideration of the uses to which differences may be put, it is instructive to consider some existing mathematical tables and their relation to differences.

As a first example consider the following section of a typical schoolboy's 4-place logarithm table:

x											Proportional Parts									
	0	1	2	3	4	5	6	7	8	9	1	2	3	4	5	6	7	8	9	
10	0000	0043	0086	0128	0170	0212	0253	0294	0334	0374	4	8	12	17	21	25	29	33	37	
11	0414	0453	0492	0531	0569	0607	0645	0682	0719	0755	4	8	11	15	19	23	26	30	34	
12	0792	0828	0864	0899	0934	etc.					etc.									

This is a typical example of what may be called a 'two dimensional' table in that, to each lattice point of a two dimensional co-ordinate system is assigned a functional value. Such tables are usually only available for the most elementary functions and at a very low level of precision. First consider the differences of the function at the finest interval available in the table, namely $\cdot 001$.

x	$f(x)$	Δ	Δ^2
100	0000	0043	0
101	0043	0043	-1
102	0086	0042	0
103	0128	0042	
104	0170		

It is intuitively clear that, since the first difference is sensibly constant, a linear relationship exists between the function and its argument in the intervals between tabulated values. It follows that, at any rate approximately, if the value of $f(x)$ at a non-tabulated value is required no great error will result from assuming:

$$f(x_0 + \delta) = f(x_0) + \delta \cdot [f(x_1) - f(x_0)]$$

or

$$f(x_0 + \delta) = f(x_0) + \delta \cdot \Delta[f(x_0)]$$

where δ is a *proportion* of the interval $(x_1 - x_0)$, and in any case the error will not exceed unity [*i.e.* the value of $\Delta^2\{f(x_0)\}$]. On the other hand, consider the table of 'proportional parts'. These purport to contain the values of $\delta\Delta[f(x_0)]$ etc. for $\delta = \cdot 0001, \cdot 0002 \dots \cdot 0009$. But since:

$$\Delta f(100) = 0043$$

$$\Delta f(109) = 0040$$

it is evident that the proportional parts can only be approximate and

that the value at $\delta = \cdot 0009$ should be 39 for $x = 1009$ and 36 for $x = 1099$. More honest table makers would mark these proportional parts to indicate that they are not accurate over the whole range of argument which they cover.

Next consider the differences over a larger interval (01).

x	$f(x)$	Δ	Δ^2	Δ^3
10	0000	0414	-0036	+0005
11	0414	0378	-0031	+0006
12	0792	0347	-0025	+0003
13	1139	0322	-0022	
14	1461	0300		
15	1761			

It is seen that the first differences are no longer constant so that an attempt to evaluate $f(103)$, say, by the 'proportional part' technique would not be justified (it would lead to $\cdot 3 \times 0414 = 0124$ which does not agree at all well with the actual value 0128), and a more sophisticated technique would be needed to make full use of the 4-place tabular accuracy.

The virtue of the simple table, described above, lies in the ease with which it may be used; when more accurate values are required, however, the two dimensional layout is no longer possible. Thus, to continue further the logarithm table to an accuracy of 7 decimal places, and to have direct reading of a 7-place argument would require a volume of some 2000 pages!

To overcome this difficulty, the onus of calculation of intermediate values is placed upon the user, and the interval of tabulation is so chosen as to make possible the linear interpolation process discussed above. (Linear interpolation is the technical term for the 'proportional part' technique just examined). An example from a 7-decimal place logarithm tabulation is the following:

x	$\log x$	Diff.
268 11	428 3130	162
12	3292	162
13	3454	162
14	3616	162
15	3778	162
16	3940	162
17	4102	162
18	4264	162
19	4426	162
20	4588	162

At a later position in the compendium from which this example is taken is to be found a table of proportional differences corresponding to all differences encountered in the main table, which in this case range from 434 to 43. A section of this table is:

Diff.	Prop.	1	2	3	4	5	6	7	8	9
162		16	32	49	65	81	97	113	130	146
163		16	33	49	65	82	98	114	130	147
164		16	33	49	66	82	98	115	131	148
165		17	33	50	66	83	99	116	132	149

Thus, to find $\log 2.68143$, the user of the table has to form the sum:

$$\begin{aligned}\log 2.68143 &= \log 2.6814 + .3 \times .0000162 \\ &= .4283616 + .0000049 \\ &= .4283665\end{aligned}$$

If greater precision is required, and careful consideration is needed to justify it in any particular case, the function difference must be actually multiplied by the proportional part. Thus:

$$\begin{aligned}\log 2.681432 &= .4283616 + .32 \times .0000162 \\ &= .4283668\end{aligned}$$

It should be noticed that when the ultimate accuracy is required, as in this case, there is a possibility of an error of ± 1 in the last place due to round off in the original table construction. This effect is eliminated in the so-called 'critical tables' in which the range of x for which $f(x)$ has a given value is specified. These tables are, however, comparatively rare and are unlikely to come the way of the student.

Our final example is taken from a recent⁽¹⁾ table of high precision (15 decimal places) for the trigonometrical function $\sin x^\circ$

x°	$\sin x$	δ^2
17.60	0.30236 98907 50445	— 92 10714
.61	.30253 62492 99766	92 15781
.62	.30270 25986 33306	92 20848
.63	.30286 89387 45998	92 25916
.64	.30303 52696 32774	92 30982
17.65	0.30320 15912 88568	— 92 36048
.66	<i>etc.</i>	

It will be noticed that, to this precision, no attempt could be made to subdivide at an interval sufficiently small to make possible linear

interpolation. Instead second central differences (see section 3.4 *infra*) have been given and these make possible a reasonable interpolation process (that of Everett, section 3.4, equation 3.4.8) for the evaluation of intermediate values.

REFERENCE

- ⁽¹⁾ Table of Sines and Cosines to 15 decimal places at hundredths of a degree.
U.S. Nat. Bur. Stand. Applied Mathematics Series. No. 5, Washington (1949)

3 INTERPOLATION

3.1 NOTATION

THE reader will be familiar with the notion of a difference as defined in section 2.1; this is not, however, the only type of difference which is convenient in numerical work. The notation:

$$\Delta u_n = u_{n+1} - u_n \quad \dots (3.1.1)$$

is usually referred to as the 'forward' difference at u_n for a reason which will be apparent from the following scheme:

x	$f(x)$	$\Delta f(x)$	$\Delta^2 f(x)$	$\Delta^3 f(x)$
0	u_0			
1	u_1	Δu_0		
2	u_2	Δu_1	$\Delta^2 u_0$	
3	u_3	Δu_2	$\Delta^2 u_1$	$\Delta^3 u_0$
4	u_4	Δu_3	$\Delta^2 u_2$	$\Delta^3 u_1$

which represents, a convenient method of tabulation.

On the other hand, it is equally possible to form the differences $(u_0 - u_{-1})$, $(u_{-1} - u_{-2})$, $(u_{-2} - u_{-3})$ etc., and it is convenient to have a 'notation' for these, although it must be borne in mind that they do not really differ from appropriate entries in the forward difference table. The customary notation for the 'backward' difference is

$$\nabla u_n = u_n - u_{n-1} \quad \dots (3.1.2)$$

and its position in tabulation is shown below:

x	$f(x)$	$\nabla f(x)$	$\nabla^2 f(x)$	$\nabla^3 f(x)$
-4	u_{-4}			
-3	u_{-3}	∇u_{-3}		
-2	u_{-2}	∇u_{-2}	$\nabla^2 u_{-2}$	
-1	u_{-1}	∇u_{-1}	$\nabla^2 u_{-1}$	$\nabla^3 u_{-1}$
0	u_0	∇u_0	$\nabla^2 u_0$	$\nabla^3 u_0$

NOTATION

At this point it should be noticed that the forward difference is particularly appropriate at the start of a table, and the backward difference at the end where, in the absence of analytical knowledge about the nature of the tabulated function, forward differences are not defined.

For intermediate points a third type of difference is appropriate, this is the so-called 'central difference' defined by :

$$\delta u_n = u_{n+\frac{1}{2}} - u_{n-\frac{1}{2}} \quad \dots (3.1.3)$$

The table then becomes:

x	$f(x)$	$\delta f(x)$	$\delta^2 f(x)$	$\delta^3 f(x)$	$\delta^4 f(x)$
-2	u_{-2}				
-1	u_{-1}	$\delta u_{-\frac{3}{2}}$			
0	u_0	$\delta u_{-\frac{1}{2}}$	$\delta^2 u_{-1}$	$\delta^3 u_{-\frac{1}{2}}$	
1	u_1	$\delta u_{\frac{1}{2}}$	$\delta^2 u_0$	$\delta^3 u_{\frac{1}{2}}$	$\delta^4 u_0$
2	u_2	$\delta u_{\frac{3}{2}}$	$\delta^2 u_1$		

which suggests the appropriateness of the central difference in the body of a table.

A further operator, which is of great utility, is defined by:

$$Eu_n = u_{n+1} \quad \dots (3.1.4)$$

or

$$Ef(x) = f(x + \delta x)$$

Successive application of 3.1.4 leads to the symbolic equation:

$$E^m u_n = u_{n+m} \quad \dots (3.1.5)$$

Again, 3.1.1 can be re-written:

$$u_{n+1} = u_n + \Delta u_n$$

whence, formally,

$$Eu_n = (1 + \Delta)u_n$$

or

$$E \equiv 1 + \Delta \quad \dots (3.1.6)$$

This relation is one of great power in deriving relationships between a function and its differences.

In a similar manner, from equation 3.1.2:

$$\nabla u_n = u_n - E^{-1}u_n$$

or

$$\nabla \equiv (E - 1)/E \quad \dots (3.1.7)$$

And, from equation 3.1.3:

$$\delta u_n = E^{\frac{1}{2}} u_n - E^{-\frac{1}{2}} u_n$$

or

$$\delta \equiv E^{\frac{1}{2}} - E^{-\frac{1}{2}} \quad \dots (3.1.8)$$

A further operator which is in common use is the 'averaging operator' μ , this is defined by:

$$\mu u_n = \frac{1}{2}(u_{n+\frac{1}{2}} + u_{n-\frac{1}{2}}) \quad \dots (3.1.9)$$

whence:

$$\mu \equiv \frac{1}{2}(E^{\frac{1}{2}} + E^{-\frac{1}{2}}) \quad \dots (3.1.10)$$

The use of the operator E is not the only method of forming $f(x + \delta x)$ from $f(x)$. Thus Taylor's theorem gives:

$$f(x + \delta x) = f(x) + \delta x f'(x) + \frac{(\delta x)^2}{2!} f''(x) + \dots$$

or, writing the differential in operator form such that:

$$f'(x) = Df(x)$$

$$f^n(x) = D^n f(x)$$

$$f(x + \delta x) = f(x) + \delta x Df(x) + \frac{(\delta x)^2}{2} D^2 f(x) + \dots$$

whence, symbolically:

$$Ef(x) = \left[1 + \delta x D + \frac{(\delta x D)^2}{2} + \dots \right] f(x)$$

or

$$E \equiv e^{\delta x D} \quad \dots (3.1.11)$$

The operator $(\delta x)D$ is sometimes represented by

$$U \equiv (\delta x)D \quad \dots (3.1.12)$$

3.2 SOME EXPANSIONS

We shall now use the notation and operators defined in section 3.1 to obtain a number of useful results. In a short account, such as the present, it is not possible to justify the operational procedures rigorously, but the reader can, if he so desires, refer to one of the standard texts, mentioned in the references, for a complete demonstration. In the meanwhile, it is worth mentioning that the operational method enables most of the standard finite-difference formulae to be worked out quickly and, as such, is a great *aide memoire* when no compilation is available.

First let us attempt to evaluate $\Delta^n u_0$ in terms of the functional values, $u_0, u_1 \dots$ etc. We have:

$$\begin{aligned} \Delta^n (u_0) &\equiv (E - 1)^n u_0 \\ &\equiv \left\{ E^n - \frac{n}{1!} E^{n-1} + \frac{n(n-1)}{2!} E^{n-2} \right. \\ &\quad \left. - \frac{n(n-1)(n-2)}{3!} E^{n-3} \dots (-1)^n \right\} u_0 \\ &= u_n - \frac{n}{1!} u_{n-1} + \frac{n(n-1)}{2!} u_{n-2} \dots (-1)^n u_0 \quad \dots (3.2.1) \end{aligned}$$

A standard notation for the r th binomial coefficient is given by:

$$n_r = \frac{n(n-1)(n-2) \dots (n-r+1)}{r!} \quad \dots (3.2.2)$$

so that equation 3.2.1 takes the simple form:

$$\Delta^n u_0 = u_n - n_1 u_{n-1} + n_2 u_{n-2} \dots (-1)^r n_r u_{n-r} \dots (-1)^n u_0 \quad \dots (3.2.3)$$

At this point it is worth pausing to discuss the application of differencing techniques to the detection of errors in tabulated functions. It is clear from equation 3.2.3 that if the m th tabulated value of a function is in error by a small quantity ϵ , the effect on the differences will increase in a regular fashion and, for the n th difference, will have a maximum effect along a central difference line pointing to the incorrect value. This is shown in the section of table below, where the presence of an error in the functional value at $x = 38.5^\circ$ is clearly indicated. The magnitude and mode of

x°	$\sin x$	δ	δ^2	δ^3	δ^4
38.0	.6157				
.1	.6170	13			
.2	.6184	14	1		
.3	.6198	14	0	-1	0
.4	.6211	13	-1	-1	-2
		9	-4	-3	17
				14	
.5	.6220		10		-30
.6	.6239	19	-6	-16	23
.7	.6252	13	1	7	-8
.8	.6266	14	0	-1	
.9	.6280	14			

INTERPOLATION

propagation of the disturbance produced by a single unit error is shown more clearly below:

f	δ	δ^2	δ^3	δ^4
0				
0	0			
0	0	0		
0	0	0	0	1
0	0	1	1	-4
1	1	-2	-3	6
0	-1	1	3	-4
0	0	0	-1	1
0	0	0	0	
0	0	0		
0	0			

A second expansion is obtainable as follows. Consider u_n , this may be written:

$$u_n = E^n u_0 = (1 + \Delta)^n u_0 \\ = u_0 + n_1 \Delta u_0 + n_2 \Delta^2 u_0 + \dots + n_r \Delta^r u_0 + \dots + \Delta^n u_0 \quad \dots (3.2.4)$$

from which it can be seen that if the first n leading differences of a function are given (for a known interval), and if the initial value is known, n tabular values are calculable. An extension of this occurs when the n th differences of a function are known to be constant. In this event all differences of order greater than n are zero, and 3.2.4 thus enables a complete tabulation to be made. Actually direct application of 3.2.4 is seldom made since it is easier to proceed directly from the table; this appears below in the tabulation of x^3 , a function for which third differences are constant.

x	x^3	Δ	Δ^2	Δ^3	Δ^4
0	0	1	6	6	0
1	1	7	12	6	
2	8	19	18	6	
3	27	37	24		
4	64	61			
5	125				
etc.					

INTERPOLATION FORMULAE

The assumption of constancy of r th differences is easily seen to imply that u is a function of degree r in the interval. This is clear from equation 3.2.4 since, when $\Delta^r u_0 = \text{const}$, $\Delta^{r+k} u_0 = 0$ ($k = 1 \dots$), and thus n_{r+k} will be multiplied by a zero coefficient. Since n_r is of degree r in the interval, the result follows. An alternative demonstration is the following:

$$\Delta x^r = [x + (\delta x)]^r - x^r \\ = r(\delta x)x^{r-1} + O(x^{r-2}) \quad \dots (3.2.5)$$

Thus the operation Δ on x^r lowers the degree of the function by unity. Repeating r times:

$$\Delta^r x^r = r! (\delta x)^r \quad \dots (3.2.6)$$

which is the required result.

A frequently used notation in finite difference work is

$$x^{(m)} = x(x-1)(x-2) \dots (x-m+1) \dots \quad (3.2.7)$$

we leave it as an exercise to the reader to prove that:

$$\Delta x^{(m)} = m x^{(m-1)} \quad \dots (3.2.8)$$

3.3 INTERPOLATION FORMULAE

We have seen that the values of a function at the points which are integral multiples of the difference interval are obtainable from equation 3.2.4. The assumption is that the function is representable by a polynomial of degree r in the interval (r th differences assumed constant), and, if it can be taken that between tabulation intervals the same polynomial is an adequate representation of the function:

$$u_x = u_0 + \xi_1 \Delta u_0 + \xi_2 \Delta^2 u_0 + \dots \quad \dots (3.3.1)$$

where

$$\xi_r = \frac{x/\delta x (x/\delta x - 1) \dots (x/\delta x - r + 1)}{r!} \quad \dots (3.3.2)$$

and δx is the interval of tabulation, so that $\xi_r = x_r$ when $\delta x = 1$.

It must be made quite clear, however, that this result, known as the 'Newton-Gregory' interpolation formula, does make the assumption of polynomial representability. To emphasize this point more strongly, assume that $r+1$ tabular values are given, so that no differences above the r th are calculable. It follows from equation 3.2.4 that

$$u_n = \sum_{s=0}^r n_s \Delta^s u_0$$

Now we may add to u_n any function:

$$\phi = n(n-1)(n-2) \dots (n-r)\psi(n)$$

where $\psi(n)$ is arbitrary, since ϕ is zero at all tabulated values and will not, in consequence, be observed. Unfortunately, however, ϕ will not, in general, be zero at non-tabulated values, e.g. at u_x ($0 < x < 1$) and any use of the Newton-Gregory formula in these circumstances will lead to error.

Even when u_n is given at *all* points ($-\infty \leq n \leq \infty$) for integral n there is still the possibility of an added function:

$$\phi = \sum_{t=0}^{\infty} a_t \sin \pi t n$$

so that care is required here also.

The reader may consider that these cases are exceptional, but the following example is one which, in modified form, might be encountered in practice, and for which Newton-Gregory interpolation is impossible.

x	$f(x)$	Δ	Δ^2	Δ^3	Δ^4
0	1	3	9	27	81
1	4	12	36	108	
2	16	48	144		
3	64	192			
4	256				

We require to find $f(\frac{1}{2})$, and the Newton-Gregory formula gives:

$$\begin{aligned} f\left(\frac{1}{2}\right) &= 1 + \frac{(\frac{1}{2})}{1!} \cdot 3 + \frac{(\frac{1}{2})(-\frac{1}{2})}{2!} \cdot 9 + \frac{(\frac{1}{2})(-\frac{1}{2})(-\frac{3}{2})}{3!} \cdot 27 \\ &\quad + \frac{(\frac{1}{2})(-\frac{1}{2})(-\frac{3}{2})(-\frac{5}{2})}{4!} \cdot 81 + \dots \\ &= 1 + \frac{3}{2} - \frac{9}{8} + \frac{27}{16} - \frac{405}{128} \dots \end{aligned}$$

from which it is evident that the series diverges and will never produce the correct value $(4)^{\frac{1}{2}} = 2$.

The Newton-Gregory formula will, in general, fail when the function cannot be expanded in a Taylor series, but in any case it is rarely useful in practice because it makes use of tabular values in *advance* of the point at which a function value is required. This is shown in the scheme:

n	f	Δ	Δ^2	Δ^3	Δ^4
-4	u_{-4}	Δu_{-4}			
-3	u_{-3}	Δu_{-3}	$\Delta^2 u_{-3}$		
-2	u_{-2}	Δu_{-2}	$\Delta^2 u_{-2}$	$\Delta^3 u_{-2}$	
-1	u_{-1}	Δu_{-1}	$\Delta^2 u_{-1}$	$\Delta^3 u_{-1}$	$\Delta^4 u_{-1}$
0	u_0	Δu_0	$\Delta^2 u_0$	$\Delta^3 u_0$	$\Delta^4 u_0$
1	u_1	Δu_1	$\Delta^2 u_1$	$\Delta^3 u_1$	$\Delta^4 u_1$
2	u_2	Δu_2	$\Delta^2 u_2$	$\Delta^3 u_2$	$\Delta^4 u_2$
3	u_3	Δu_3	$\Delta^2 u_3$	$\Delta^3 u_3$	$\Delta^4 u_3$
4	u_4	Δu_4	$\Delta^2 u_4$	$\Delta^3 u_4$	$\Delta^4 u_4$

Gauss 'backward'

Gauss 'forward'

Newton-Gregory

A more reasonable and convergent procedure would be to make use of functional values which straddle the point at which interpolation is required. Two such schemes are shown in the diagram, and these lead to the Gaussian interpolation formulae:

$$u_x = u_0 + x\Delta u_0 + x_2\Delta^2 u_{-1} + (x+1)_3\Delta^3 u_{-1} + (x+1)_4\Delta^4 u_{-2} + \dots \quad (3.3.3)$$

which is the Gauss 'forward' formula, and:

$$u_x = u_0 + x\Delta u_{-1} + (x+1)_2\Delta^2 u_{-1} + (x+1)_3\Delta^3 u_{-2} + (x+2)_4\Delta^4 u_{-2} + \dots \quad (3.3.4)$$

which is the Gauss 'backward' formula. x_r has been defined in equation 3.3.2. The formulae are also correct when x is the proportional part of the tabulation interval.

It has been remarked by COMRIE⁽¹⁾ that these formulae, and also that of Stirling are never used in good modern practice, and for this reason we do not give a detailed proof. The reader will satisfy himself as to the correctness of the first few terms in the 'forward' formula by substituting from:

$$\begin{aligned} \Delta^2 u_0 &= \Delta^2 u_{-1} + \Delta^3 u_{-1} \\ \Delta^3 u_0 &= \Delta^3 u_{-1} + \Delta^4 u_{-1} \\ \Delta^4 u_{-1} &= \Delta^4 u_{-2} + \Delta^5 u_{-2} \quad \text{etc.} \end{aligned}$$

in the Newton-Gregory result (3.3.1) and rearranging. The 'backward' formula is likewise obtained by substituting from:

$$\begin{aligned} \Delta u_0 &= \Delta u_{-1} + \Delta^2 u_{-1} \\ \Delta^2 u_0 &= \Delta^2 u_{-1} + \Delta^3 u_{-1} \\ \Delta^3 u_{-1} &= \Delta^3 u_{-2} + \Delta^4 u_{-2} \quad \text{etc.} \end{aligned}$$

A general derivation has been given by FREEMAN⁽²⁾.

If we take the mean of the two Gauss formulae (3.3.3 and 3.3.4) we obtain:

$$u_x = u_0 + \frac{1}{2}x(\Delta u_0 + \Delta u_{-1}) + \frac{x^2}{2!}\Delta^2 u_{-1} + \frac{x(x^2-1^2)}{2 \cdot 3!}(\Delta^3 u_{-1} + \Delta^3 u_{-2}) + \frac{x^2(x^2-1^2)}{4!}\Delta^4 u_{-2} + \frac{x(x^2-1^2)(x^2-2^2)}{2 \cdot 5!}(\Delta^5 u_{-2} + \Delta^5 u_{-3}) + \dots \quad (3.3.5)$$

which is known as Stirling's formula.

We now proceed to derive the two formulae which are of most common use in actual problems. Since they appear at this point they will be expressed in terms of the forward difference operator Δ . It is, however, more usual to work in central differences and an alternative proof in terms of these operators will be given later, in section 3.4.

First we transform the origin of the Gauss backward formula from u_0 to u_1 so that x becomes $(x-1)$ and:

$$u_x = u_1 + (x-1)\Delta u_0 + \frac{x(x-1)}{2!}\Delta^2 u_0 + \frac{x(x-1)(x-2)}{3!}\Delta^3 u_{-1} + \frac{(x+1)x(x-1)(x-2)}{4!}\Delta^4 u_{-1} + \dots \quad (3.3.6)$$

Taking the mean of equations 3.3.6 and 3.3.3,

$$u_x = \frac{1}{2}(u_1 + u_0) + (x - \frac{1}{2})\Delta u_0 + \frac{x(x-1)}{2 \cdot 2!}(\Delta^2 u_{-1} + \Delta^2 u_0) + \dots \quad (3.3.7)$$

which is Bessel's formula.

Finally,⁽³⁾ we may write the Gauss forward formula (equation 3.3.3):

$$w_{1+x} = w_1 + x\Delta w_1 + x_2\Delta^2 w_0 + (x+1)_3\Delta^3 w_0 + (x+1)_4\Delta^4 w_{-1} + \dots$$

and similarly equation 3.3.6:

$$w_x = w_1 + (x-1)\Delta w_0 + x_2\Delta^2 w_0 + x_3\Delta^3 w_{-1} + (x+1)_4\Delta^4 w_{-1}$$

Now if we subtract w_x from w_{1+x} we obtain:

$$\Delta w_x = x\Delta w_1 + (x+1)_3\Delta^3 w_0 + \dots - (x-1)\Delta w_0 - x_3\Delta^3 w_{-1} - \dots$$

(since $w_{1+x} - w_x = \Delta w_x$ by definition)

whence, on placing $u_x = \Delta w_x$:

$$u_x = xu_1 + (x+1)_3\Delta^2 u_0 + \dots - (x-1)u_0 - x_3\Delta^2 u_{-1} \dots$$

or

$$u_x = xu_1 + \frac{x(x^2-1)}{3!}\Delta^2 u_0 + \dots - (x-1)u_0 - \frac{(x-1)(\{x-1\}^2-1)}{3!}\Delta^2 u_{-1} \dots \quad (3.3.8)$$

which is Everett's formula.

3.4 CENTRAL DIFFERENCES

It is clear from the analysis of section 3.3 that a difference interpolation formula, for use in the body of a table, suggests the notion of central differences (which were mentioned briefly in section 3.1). We shall now examine an operational method for obtaining central difference formulae in a natural manner.

First observe, from the central difference notation table of section 3.1, that the only central differences of u_0, u_1 which are available from the tabular values are the even ones δ^2, δ^4 etc. Let us, therefore, attempt to find an expansion for u_x in terms of $\delta^{2n}u_0$ and $\delta^{2n}u_1$.

Operationally we have,

$$u_x = E^x u_0 = e^{xU} \cdot u_0 \quad (3.4.1)$$

where U is defined by equation 3.1.12.

Assuming that an expansion is possible in terms of $\delta^{2n}u_0$ and $\delta^{2n}u_1$, this will be of the form:

$$F(\delta) \cdot u_0 + G(\delta)u_1 \quad (3.4.2)$$

where F and G are even functions of δ .

Now

$$u_1 = Eu_0 = e^U u_0 \quad (3.4.3)$$

whence, from equations 3.4.1, 3.4.2 and 3.4.3:

$$e^{xU} \equiv F(\delta) + G(\delta)e^U \quad (3.4.4)$$

Now from equations 3.1.8 and 3.1.11 we have:

$$\frac{\delta}{2} = \sinh \left(\frac{U}{2} \right) \quad (3.4.5)$$

so that δ is an odd function of U . Now, since F and G are even

functions of δ , they are also even functions of U . Thus, replacing U by its negative in equation 3.4.4, we obtain:

$$e^{-xU} = F(\delta) + G(\delta)e^{-U} \quad \dots (3.4.6)$$

and, after some manipulation:

$$\left. \begin{aligned} F(\delta) &= \sinh (1-x)U / \sinh U \\ G(\delta) &= \sinh xU / \sinh U \end{aligned} \right\} \quad \dots (3.4.7)$$

From these expressions, together with equation 3.4.5 it can be shown that:

$$F(\delta) = (1-x) \left\{ 1 + \frac{[(1-x)^2 - 1^2]}{3!} \delta^2 + \frac{[(1-x)^2 - 1^2][(1-x)^2 - 2^2]}{5!} \delta^4 + \dots \right\}$$

and that

$$G(\delta) = x \left\{ 1 + \frac{(x^2 - 1^2)}{3!} \delta^2 + \frac{(x^2 - 1^2)(x^2 - 2^2)}{5!} \delta^4 + \dots \right\}$$

so that:

$$\begin{aligned} u_x &= (1-x) \left\{ u_0 + \frac{[(1-x)^2 - 1^2]}{3!} \delta^2 u_0 + \frac{[(1-x)^2 - 1^2][(1-x)^2 - 2^2]}{5!} \delta^4 u_0 + \dots \right\} \\ &+ x \left\{ u_1 + \frac{(x^2 - 1^2)}{3!} \delta^2 u_1 + \frac{(x^2 - 1^2)(x^2 - 2^2)}{5!} \delta^4 u_1 + \dots \right\} \end{aligned} \quad \dots (3.4.8)$$

which is Everett's formula in central difference form. The reader will see, by renaming the forward difference symbols in equation 3.3.8 that it is identical with equation 3.4.8.

Bessel's formula can also be obtained in central difference form by means of the operational equation:

$$E^x u_0 = F(\delta) (u_1 + u_0) + G(\delta) u_1$$

the form of which is suggested by equation 3.3.7, and where F is again an *even* function of δ but $G(\delta)$ is this time *odd*. The final result is:

$$\begin{aligned} u_x &= \frac{1}{2}(u_0 + u_1) + (x - \frac{1}{2})\delta u_{\frac{1}{2}} + B^{II}(x)(\delta^2 u_0 + \delta^2 u_1) + B^{III}(x)\delta^3 u_{\frac{1}{2}} \\ &+ B^{IV}(x)(\delta^4 u_0 + \delta^4 u_1) + B^V(x)\delta^5 u_{\frac{1}{2}} + \dots \end{aligned} \quad \dots (3.4.9)$$

where the symbols $B^n(x)$ are given by:

$$\left. \begin{aligned} B^{II}(x) &= x(x-1)/2.2! \\ B^{III}(x) &= x(x-\frac{1}{2})(x-1)/3! \\ B^{IV}(x) &= (x+1)x(x-1)(x-2)/2.4! \\ B^V(x) &= (x+1)x(x-\frac{1}{2})(x-1)(x-2)/5! \\ &\quad \text{etc.} \end{aligned} \right\} \quad \dots (3.4.10)$$

The notation used here is that suggested by Comrie, *loc. cit.*

3.5 MODIFIED DIFFERENCES

In use, Bessel's formula is often modified by means of a numerical accident. Noting that $\delta u_{\frac{1}{2}} = u_1 - u_0$, we see that equation 3.4.9 can be written, to fourth differences:

$$\begin{aligned} u_x &= u_0 + x\delta u_{\frac{1}{2}} + B^{II}(x)(\delta^2 u_0 + \delta^2 u_1) \\ &+ B^{III}(x)\delta^3 u_{\frac{1}{2}} + B^{IV}(x)(\delta^4 u_0 + \delta^4 u_1) \end{aligned} \quad \dots (3.5.1)$$

Now reference to equation 3.4.10 shows that:

$$B^{IV}(x) = \frac{(x+1)(x-2)}{12} B^{II}(x) \quad \dots (3.5.2)$$

and a calculation shows that the quadratic coefficient of $B^{II}(x)$ has only a limited variation in the range $(0 \leq x \leq 1)$. This is shown in Figure 3.5.1.

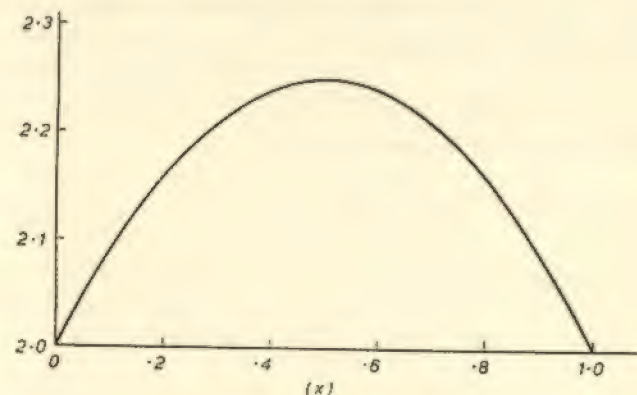


Figure 3.5.1.—12 $B^{IV}(x)/B^{II}(x)$

The limits of variation are between -2.0 at $x=0$, $x=1$ and -2.25 at $x=\frac{1}{2}$. It follows that

$$B^{iv}(x) = CB^{ii}(x)$$

where C varies between $-1.66\bar{6}$ and -1.875 .

It was suggested by L. J. Comrie that in many cases it would be adequate to write:

$$B^{iv}(x) = CB^{ii}(x)$$

where C is a constant known as the throwback coefficient, so chosen as to make the greatest absolute error, E ,

$$\text{in } E = B^{iv}(x) - CB^{ii}(x)$$

as small as possible. It is fairly easy to see that the required value of C is one which makes $E_{\max} = -E_{\min}$ in the range $0 \leq x \leq 1$. To determine the value of C , we have from 3.4.10:

$$E = \frac{x(x-1)}{4 \cdot 12} [(x+1)(x-2) - 12C]$$

We notice that E is symmetrical about $x=\frac{1}{2}$ and put $z = x - \frac{1}{2}$ so that

$$\begin{aligned} E &= \frac{(z + \frac{1}{2})(z - \frac{1}{2})}{48} [(z + \frac{3}{2})(z - \frac{3}{2}) - 12C] \\ &= \frac{(z^2 - \frac{1}{4})}{48} [z^2 - \frac{9}{4} - 12C] \end{aligned}$$

Evidently one of the maximum absolute values of E occurs at $z=0$ when $E = \frac{\frac{9}{4} + 12C}{4 \times 48}$. To determine the other maximum absolute deviation we put $y = z^2$ so that:

$$E = \frac{1}{48} [y^2 - (\frac{9}{2} + 12C)y + (\frac{9}{16} + 3C)]$$

the greatest value of E is given by $dE/dy = 0$

$$\text{i.e. } y = \frac{9}{4} + 6C$$

$$\text{whence } E_{\max} = -\frac{1}{48}(6C + 1)^2$$

We thus require

$$\frac{1}{48}(6C + 1)^2 = \frac{1}{4 \times 48} (\frac{9}{4} + 12C)$$

$$\text{or } 36C^2 + 9C + \frac{7}{16} = 0$$

whence

$$C = -3 \pm \sqrt{2}/24$$

To determine the appropriate sign we notice that $C = -3 + \sqrt{2}/24$ gives a larger absolute value of E at $z=0$ than does $C = -3 - \sqrt{2}/24$ so that the appropriate value is $C = -3 - \sqrt{2}/24$ or $C \approx -1.84$.

Using this result, fourth differences which do not exceed 1000 units in the last decimal place are 'thrown back' into the second difference to give a modified second difference defined by:

$$\delta_m^2 u_n = \delta^2 u_n - .184 \delta^4 u_n \quad \dots (3.5.3)$$

The Bessel formula is then simply written

$$u_x = u_0 + x \delta u_{\frac{1}{2}} + B^{ii}(x) (\delta_m^2 u_0 + \delta_m^2 u_1) + B^{iii}(x) \delta^3 u_{\frac{1}{2}}$$

and is then accurate to fourth central differences.

This idea of throwback can be extended to higher differences than the fourth but the algebra for determining the throwback coefficients becomes prohibitive. KOPAL⁽⁴⁾ has given a general method for deriving throwback coefficients which depend upon Chebyshev polynomials, but it is debatable if the results are superior to those which can be obtained by using a least squares approach in which C is determined by making

$$\int_0^1 E^2 dx$$

a minimum. This technique is discussed in Chapter 10.

3.6 DIVIDED DIFFERENCES

When the values of a function are given for values of the argument which are not equally spaced, it is still possible to define a system of differences. Consider the table:

x	$f(x)$
x_0	u_0
x_1	u_1
x_2	u_2
x_3	u_3

The first divided differences are then defined to be:

$$\begin{aligned} \Delta' u_0 &= (u_1 - u_0)/(x_1 - x_0) \\ \Delta' u_1 &= (u_2 - u_1)/(x_2 - x_1) \\ &\dots \\ \Delta' u_r &= (u_{r+1} - u_r)/(x_{r+1} - x_r) \quad \dots (3.6.1) \end{aligned}$$

Similarly, for second divided differences,

$$\Delta'^2 u_0 = (\Delta' u_1 - \Delta' u_0) / (x_2 - x_0)$$

$$\Delta'^2 u_1 = (\Delta' u_2 - \Delta' u_1) / (x_3 - x_1)$$

etc.

and, in general:

$$\Delta'^{r+1} u_0 = (\Delta'^r u_1 - \Delta'^r u_0) / (x_{r+1} - x_0)$$

$$\Delta'^{r+1} u_n = (\Delta'^r u_{n+1} - \Delta'^r u_n) / (x_{n+r+1} - x_n) \quad \dots (3.6.2)$$

There exists the analogue of the Newton-Gregory formula for divided differences:

$$u_x = u_0 + (x - x_0)\Delta' u_0 + (x - x_0)(x - x_1)\Delta'^2 u_0 + (x - x_0)(x - x_1)(x - x_2)\Delta'^3 u_0 + \dots \\ (x - x_0)(x - x_1) \dots (x - x_{r-1})\Delta'^r u_0 + \dots (3.6.3)$$

It can be seen from equations 3.6.1 and 3.6.2 that when the intervals $x_{r+1} - x_r$ become equal (to δx , say) the divided difference formulae reduce to:

$$\Delta' u_n = \Delta u_n / \delta x$$

$$\Delta'^2 u_n = \Delta^2 u_n / 1.2.(\delta x)^2$$

$$\Delta'^3 u_n = \Delta^3 u_n / 1.2.3(\delta x)^3$$

etc.

$$\Delta'^r u_n = \Delta^r u_n / r!(\delta x)^r \quad \dots (3.6.4)$$

so that equation 3.6.3. may be written :

$$u_x = u_0 + \frac{x\Delta u_0}{\delta x} + x(x - \delta x) \frac{\Delta^2 u_0}{2!(\delta x)^2} + x(x - \delta x)(x - 2\delta x) \frac{\Delta^3 u_0}{3!(\delta x)^3} + \dots \\ + x(x - \delta x)(x - 2\delta x) \dots [x - (r - 1)\delta x] \frac{\Delta^r u_0}{r!(\delta x)^r} + \dots$$

or

$$u_x = u_0 + \xi_1 \Delta u_0 + \xi_2 \Delta^2 u_0 + \dots \xi_r \Delta^r u_0 + \dots$$

where ξ_r is as defined in equation 3.3.2.

This is the ordinary Newton-Gregory formula of equation 3.3.1 for interval δx .

3.7 LAGRANGEAN INTERPOLATION

An alternative formula for interpolation at unequal intervals is due to Lagrange. Suppose that, as before, functional values:

x	$f(x)$
x_0	u_0
x_1	u_1
x_2	u_2
\vdots	\vdots
x_{n-1}	u_{n-1}

are given, and that they are n in number. An interpolation formula consists in finding a polynomial, of degree $(n - 1)$, passing through the given values. Assume that:

$$f(x) = (x - x_0)(x - x_1) \dots (x - x_{n-1}) \sum_{r=0}^{n-1} \frac{A_r}{(x - x_r)} \quad \dots (3.7.1)$$

Put $x = x_s$ ($s = 0, 1 \dots n - 1$) then:

$$A_s = \frac{f(x_s)}{(x_s - x_0)(x_s - x_1) \dots (x_s - x_{s-1})(x_s - x_{s+1}) \dots (x_s - x_{n-1})} \\ = u_s / (x_s - x_0)(x_s - x_1) \dots (x_s - x_{s-1})(x_s - x_{s+1}) \dots (x_s - x_{n-1}) \quad \dots (3.7.2)$$

Equations 3.7.1 and 3.7.2, together, constitute Lagrange's interpolation formula.

3.8 CAUTIONS AND PRECAUTIONS

One example has already been given (section 3.3) of the dangers which arise from the blind use of direct interpolation formulae. The reverse process, called inverse interpolation, and sufficiently defined by the order: 'given $f(x)$, find x ', is perhaps even more strewn with pitfalls for the unwary. As an indication of these, suppose that the values:

x	$f(x)$
1	0
4	1
16	2
64	3
256	4

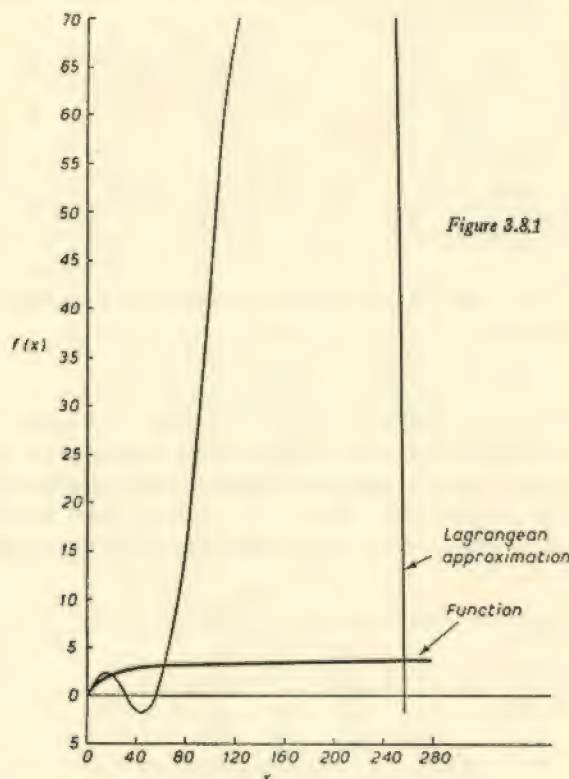
are given, and that it is desired to find $f(x)$ when $x = 32$. (The reader will notice that this is really inverse interpolation from the table given previously in section 3.3.)

The appropriate Lagrange formula is:

INTERPOLATION

$$f(x) = 1 \cdot \frac{(x-1)(x-16)(x-64)(x-256)}{(4-1)(4-16)(4-64)(4-256)} \\ + 2 \cdot \frac{(x-1)(x-4)(x-64)(x-256)}{(16-1)(16-4)(16-64)(16-256)} \\ + 3 \cdot \frac{(x-1)(x-4)(x-16)(x-256)}{(64-1)(64-4)(64-16)(64-256)} \\ + 4 \cdot \frac{(x-1)(x-4)(x-16)(x-64)}{(256-1)(256-4)(256-16)(256-64)}$$

from which it is easily verified that $f(32) = -0.263$, which differs considerably from the true value $+2.5$!



The reason for this behaviour is quite evident when the mode of variation of the function, and of its Lagrangean approximation, are considered. The two expressions are shown in Figure 3.8.1.

CAUTIONS AND PRECAUTIONS

In this particular example it happens that a linear interpolation, in the range concerned, would give a far better result. This emphasizes the fact that, before using an interpolation formula on any tabulated function, it should be carefully ascertained that the function behaves in a manner capable of representation by the formula.

REFERENCES

- (1) COMRIE, L. J., Chambers's Six-Figure Tables, 2 (1949) p. xxviii
- (2) FREEMAN, H., *J. Inst. Actuaries*, L, (1924), 31
- (3) LIDSTONE, J., *J. Inst. Actuaries*, LX (1934) 349
- (4) KOPAL, Z., 'Numerical Analysis,' p. 54, Chapman & Hall, London (1955)

NUMERICAL DIFFERENTIATION AND INTEGRATION

4.1 OPERATORS

The differential of a function has already been introduced as D in section 3.1. In a similar manner we may introduce an operational symbol for integration as:

$$\int \equiv \frac{1}{D} (\equiv D^{-1}) \quad \dots (4.1.1)$$

Now from equation 3.1.11 we have:

$$\delta x D \equiv \log_e E = \log_e (1 + \Delta).$$

whence, expanding:

$$\delta x D \equiv \Delta - \frac{\Delta^2}{2} + \frac{\Delta^3}{3} - \dots$$

or

$$D \equiv \frac{\Delta}{\delta x} \left[1 - \frac{\Delta}{2} + \frac{\Delta^2}{3} \dots + (-1)^r \frac{\Delta^r}{r+1} + \dots \right] \dots (4.1.2)$$

which, when applied to $y = f(x)$ at $x = x_0$ gives:

$$\left(\frac{dy}{dx} \right)_{x=x_0} = \frac{1}{\delta x} [\Delta f(x_0) - \frac{1}{2} \Delta^2 f(x_0) + \frac{1}{3} \Delta^3 f(x_0) - \dots] \quad \dots (4.1.3)$$

In a like manner, raising 4.1.2 to power n .

$$D^n \equiv \frac{\Delta^n}{(\delta x)^n} \left(1 - \frac{\Delta}{2} + \frac{\Delta^2}{3} \dots \right)^n \quad \dots (4.1.4)$$

from which:

$$D^2 \equiv \frac{\Delta^2}{(\delta x)^2} \left(1 - \Delta + \frac{11}{12} \Delta^2 - \frac{5}{6} \Delta^3 + \dots \right) \quad \dots (4.1.5)$$

and

$$D^3 \equiv \frac{\Delta^3}{(\delta x)^3} \left(1 - \frac{3}{2} \Delta + \frac{7}{4} \Delta^2 \dots \right) \quad \dots (4.1.6)$$

These formulae provide a convenient means of obtaining the differential coefficients of a function which is given by a table of values, and not in analytical form.

4.2 CENTRAL DIFFERENCE FORMULAE FOR DIFFERENTIATION

In many applications it is more convergent to have expressions for the differential coefficients of a function in terms of central differences. The operational method may again be applied in the following manner.

From equation 3.4.5:

$$U(\equiv \delta x D) \equiv 2 \sinh^{-1} \frac{\delta}{2}$$

where δ is the central difference operator.

Unfortunately U is an odd function of δ so that this equation will produce a series of *odd* powers of δ for all *odd* derivatives. Since such values of δ cannot be obtained directly from tabular entries (see the third diagram of section 3.1) a more satisfactory procedure has to be found. This can be done by means of the operator μ , defined by equation 3.1.10.

Clearly

$$\mu \delta \equiv \frac{1}{2}(E^1 - E^{-1}) \quad \dots (4.2.1)$$

so that

$$\mu \delta u_0 = \frac{1}{2}(u_1 - u_{-1})$$

which can be found directly from a table.

Also:

$$2\mu + \delta \equiv 2E^{\frac{1}{2}} \quad \dots (4.2.2)$$

$$2\mu - \delta \equiv 2E^{-\frac{1}{2}} \quad \dots (4.2.3)$$

whence

$$4\mu^2 - \delta^2 \equiv 4$$

or

$$\mu^2 \equiv \frac{1}{4}\delta^2 + 1 \quad \dots (4.2.5)$$

Now we may write:

$$\left(\frac{U}{\delta} \right)^n \equiv \left(\frac{\sinh^{-1} \frac{1}{2}\delta}{\frac{1}{2}\delta} \right)^n \quad \dots (4.2.6)$$

which will give, for *even* values of n , an expansion for D^n in terms of *even* powers of δ . Alternatively:

$$\frac{1}{\mu} \left(\frac{U}{\delta} \right)^n = \frac{1}{\mu} \left(\frac{\sinh^{-1} \frac{1}{2}\delta}{\frac{1}{2}\delta} \right)^n \quad \dots (4.2.7)$$

which, by virtue of equation 4.2.5, will give expansions of D^n for *odd* values of n , in terms of *even* powers of δ multiplied by $\mu \delta^n$ which, since n is *odd*, can be found from tabular entries.

To obtain the series we make use of two general results due to BICKLEY⁽¹⁾ :

$$\begin{aligned} \left(\frac{U}{\delta}\right)^n &= 1 - \frac{n}{24}\delta^2 + \frac{5n^2 + 22n}{5760}\delta^4 - \frac{35n^3 + 462n^2 + 1528n}{2903040}\delta^6 + \\ &+ \frac{175n^4 + 4620n^3 + 40724n^2 + 119856n}{1393459200}\delta^8 - \\ &- \frac{385n^5 + 16940n^4 + 279884n^3 + 2057968n^2 + 5682048n}{367873228800}\delta^{10} \\ &+ \dots \quad \dots (4.2.8) \end{aligned}$$

$$\begin{aligned} \frac{1}{\mu}\left(\frac{U}{\delta}\right)^n &= 1 - \frac{n+3}{24}\delta^2 + \frac{5n^2 + 52n + 135}{5760}\delta^4 - \\ &- \frac{35n^3 + 777n^2 + 5749n + 14175}{2903040}\delta^6 + \\ &+ \frac{175n^4 + 6720n^3 + 96794n^2 + 619776n + 1488375}{1393459200}\delta^8 - \\ &- \frac{385n^5 + 22715n^4 + 536294n^3 + 6333250n^2 + 37408281n + 88409475}{367873228800}\delta^{10} \\ &+ \dots \quad \dots (4.2.9) \end{aligned}$$

These results hold for positive and negative values of n .

From equation 4.2.8 we have:

$$\left(\frac{U}{\delta}\right)^2 \equiv 1 - \frac{1}{12}\delta^2 + \frac{1}{90}\delta^4 - \dots$$

which gives:

$$D^2 \equiv \frac{1}{(\delta x)^2} (\delta^2 - \frac{1}{12}\delta^4 + \frac{1}{90}\delta^6 - \dots) \quad \dots (4.2.10)$$

Also from equation 4.2.9:

$$\frac{1}{\mu}\left(\frac{U}{\delta}\right) \equiv 1 - \frac{1}{6}\delta^2 + \frac{1}{90}\delta^4 - \frac{1}{140}\delta^6 + \dots$$

leading to:

$$D \equiv \frac{1}{(\delta x)} [(\mu\delta) - \frac{1}{6}\delta^2(\mu\delta) + \frac{1}{90}\delta^4(\mu\delta) - \frac{1}{140}\delta^6(\mu\delta) + \dots] \quad \dots (4.2.11)$$

When applied to a functional value (u_0 , say) 4.2.11 may be written (by virtue of equation 4.2.1)

$$\begin{aligned} D(u_0) &= \frac{1}{2(\delta x)} [(u_1 - u_{-1}) - \frac{1}{6}(\delta^2 u_1 - \delta^2 u_{-1}) \\ &+ \frac{1}{90}(\delta^4 u_1 - \delta^4 u_{-1}) - \dots] \quad \dots (4.2.12) \end{aligned}$$

which is in a suitable form for tabular use.

Higher derivatives can be evaluated in a similar manner by the use of equations 4.2.8 and 4.2.9.

4.3 NUMERICAL INTEGRATION

In many problems it becomes necessary to evaluate the integral of a given function between certain limits, and although the derivatives can always be found when the analytic form of the function is given, the same is not necessarily (or even frequently) true of the integral.

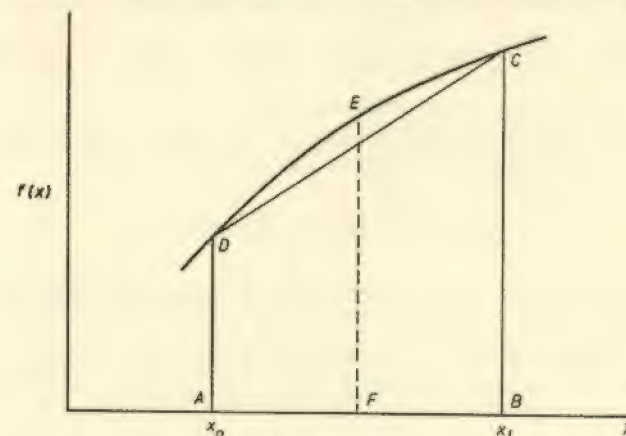


Figure 4.3.1

Thus numerical integration, sometimes called 'mechanical quadrature', is a process of importance in its own right, and quite apart from its use with tabulated functions.

The simplest form of numerical integration may be described as trapezoidal approximation. Thus, with reference to Figure 4.3.1, it is evident that a rough approximation to the integral:

$$\int_{x_0}^{x_1} f(x) dx$$

is given by the area $ABCD$, i.e. by:

$$\begin{aligned} & \frac{1}{2}(AD + BC)(x_1 - x_0) \\ &= \frac{1}{2}(x_1 - x_0)(f(x_1) + f(x_0)) \end{aligned}$$

Clearly, a better approximation may be obtained by representing the arc DEC by means of a quadratic expression in x , and forming the integral for the resulting function.

The quadratic function can easily be obtained from the Newton-Gregory formula, and is:

$$f_q(x) = f(0) + x\Delta f(0) + \frac{x(x-1)}{2}\Delta^2 f(0) \quad \dots (4.3.1)$$

where $f_q(x)$ is the quadratic approximation to $f(x)$, and we have taken the intervals $AF = FB = 1$. Equation 4.3.1 can be integrated to give:

$$\int_0^2 f_q(x) dx = 2f(0) + 2\Delta f(0) + \frac{1}{3}\Delta^2 f(0) \quad \dots (4.3.2)$$

or, replacing the differences by their tabular values:

$$\int_0^2 f_q(x) dx = \frac{1}{3}[f(2) + 4f(1) + f(0)] \quad \dots (4.3.3)$$

which is usually known as 'Simpson's rule'.

The extension to the case in which the interval in x is (δx) instead of unity is immediate and leads to the formula:

$$\int_0^{2(\delta x)} f_q(x) dx = \frac{(\delta x)}{3} [f(2) + 4f(1) + f(0)] \quad \dots (4.3.4)$$

When a function is to be integrated over an extensive range it is, in general, better to add the contributions arising from the application of Simpson's rule to successive groups of three ordinates, rather than to seek an approximating polynomial which represents the function over the whole range (see, for example section 3.8). Thus:

$$\begin{aligned} \int_0^{2n(\delta x)} f(x) dx &\approx \int_0^{2(\delta x)} f_{q_1}(x) dx + \int_{2(\delta x)}^{4(\delta x)} f_{q_2}(x) dx + \dots \\ &\quad + \int_{(2n-2)(\delta x)}^{2n(\delta x)} f_{q_n}(x) dx \\ &= \frac{(\delta x)}{3} [f(0) + 4f(1) + 2f(2) + 4f(3) + \dots + f(2n)] \quad \dots (4.3.5) \end{aligned}$$

from equation 4.3.4. This expression is often referred to as the *extension* of Simpson's rule. It should be noted that the approximating function in this case, is not necessarily a smooth curve.

A further approximation to the integral of a function, between limits 0 and 3 can be obtained by adding a further term to the Newton-Gregory formula (4.3.1). We then obtain:

$$f_c(x) = f(0) + x\Delta f(0) + \frac{x(x-1)}{2}\Delta^2 f(0) + \frac{1}{6}x(x-1)(x-2)\Delta^3 f(0)$$

where $f_c(x)$ is a *cubic* approximation to $f(x)$. Integrating we obtain:

$$\int_0^3 f_c(x) dx = 3f(0) + \frac{3}{2}\Delta f(0) + \frac{3}{4}\Delta^2 f(0) + \frac{3}{8}\Delta^3 f(0) \quad \dots (4.3.6)$$

or, writing the values of Δ , Δ^2 and Δ^3 in tabular form:

$$\int_0^3 f_c(x) dx = \frac{3}{8}[f(0) + 3f(1) + 3f(2) + f(3)] \quad \dots (4.3.7)$$

which is called the 'three eighths' rule.

A more accurate numerical integration procedure is 'Weddle's rule', which may be derived as follows. Take the first 6 differences and apply the Newton-Gregory expansion:

$$\begin{aligned} f_s(x) &= f(0) + x\Delta f(0) + \frac{x(x-1)}{2!}\Delta^2 f(0) + \frac{x(x-1)(x-2)}{3!}\Delta^3 f(0) \\ &\quad + \frac{x(x-1)(x-2)(x-3)}{4!}\Delta^4 f(0) + \frac{x(x-1)\dots(x-4)}{5!}\Delta^5 f(0) \\ &\quad + \frac{x(x-1)\dots(x-5)}{6!}\Delta^6 f(0) \end{aligned}$$

Integrating this expression over the range 0 — 6 we obtain:

$$\begin{aligned} \int_0^6 f_s(x) dx &= 6f(0) + 18\Delta f(0) + 27\Delta^2 f(0) + 24\Delta^3 f(0) \\ &\quad + \frac{123}{10}\Delta^4 f(0) + \frac{33}{10}\Delta^5 f(0) + \frac{41}{140}\Delta^6 f(0) \quad \dots (4.3.8) \end{aligned}$$

Now, with an error of $\Delta^6/140$, we may write the last term in equation 4.3.8 as $\frac{41}{140}\Delta^6 f(0) = \frac{3}{10}\Delta^6 f(0)$. The expression then becomes, on substituting for the differences in terms of the tabular values:

$$\int_0^6 f_s(x) dx = \frac{3}{10}(u_0 + 5u_1 + u_2 + 6u_3 + u_4 + 5u_5 + u_6) \quad \dots (4.3.9)$$

which is Weddle's formula in the usual notation.

The foregoing results are sometimes known as Newton-Cotes integration formulae of the *closed* type. The expression *closed* is inserted to indicate that the approximating polynomials pass through the end points or limits of integration.

For the purpose of solving differential equations it is useful to have available integration formulae of the Newton-Cotes *open* type. That corresponding to Simpson's rule can be obtained by integrating the quadratic approximation 4.3.1, shifted to operate on $x = 1$ as base point. Thus:

$$f_q(x) = f(1) + (x-1) \Delta f(1) + \frac{(x-1)(x-2)}{2} \Delta^2 f(1)$$

and

$$\int_0^4 f_q(x) dx = \frac{4}{3} [2f(1) - f(2) + 2f(3)]$$

or, for interval (δx)

$$\int_0^{4(\delta x)} f_q(x) dx = \frac{4}{3} (\delta x) [2f(1) - f(2) + 2f(3)] \dots (4.3.10)$$

4.4 THE EULER-MACLAURIN EXPANSION

A general integration formula which is often quoted in numerical analysis is that due to Euler and Maclaurin.

Consider the summation:

$$\sum_{x=a}^{x=a+(n-1)\delta x} f(x) = \left(\sum_{r=0}^{n-1} E^r \right) f(a) \dots (4.4.1)$$

Proceeding in terms of operators we have:

$$\sum_{r=0}^{n-1} E^r \equiv \frac{E^n - 1}{E - 1} \dots (4.4.2)$$

Now, by virtue of equation 3.1.11 $E = e^{(\delta x D)} \equiv e^U$ and we may write:

$$\frac{1}{E - 1} \equiv \frac{1}{e^U - 1} \equiv \frac{1}{U} \left[\frac{U}{e^U - 1} \right] \dots (4.4.3)$$

and the well-known ⁽²⁾ expansion

$$\frac{U}{e^U - 1} = 1 - \frac{U}{2} + B_1 \frac{U^2}{2!} - B_2 \frac{U^4}{4!} + B_3 \frac{U^6}{6!} - \dots$$

$$+ (-1)^{r+1} B_r \frac{U^{2r}}{(2r)!} + \dots$$

where B_r is the r th Bernoulli number, enables equation 4.4.3 to be written:

$$\frac{1}{E - 1} \equiv \frac{1}{U} - \frac{1}{2} + B_1 \frac{U}{2!} - B_2 \frac{U^3}{4!} + B_3 \frac{U^5}{6!} - \dots (4.4.4)$$

The numerical values of the first few Bernoulli numbers are:

$$B_1 = \frac{1}{6}, \quad B_2 = \frac{1}{30}, \quad B_3 = \frac{1}{42}, \quad B_4 = \frac{1}{30}, \quad B_5 = \frac{5}{66},$$

$$B_6 = \frac{691}{2730}, \quad B_7 = \frac{7}{6}, \quad B_8 = \frac{3817}{510}, \quad B_9 = \frac{43847}{798}, \quad B_{10} = \frac{174611}{330}$$

Substituting from equations 4.4.4 and 4.4.2 in equation 4.4.1, and

observing that $U = (\delta x D)$, $\frac{1}{U} = \frac{1}{\delta x} \int$, we obtain:

$$\sum_{x=a}^{x=a+(n-1)\delta x} f(x) =$$

$$= (E^n - 1) \left[\frac{1}{\delta x} \int f(a) da - \frac{1}{2} f(a) + \frac{B_1}{2!} (\delta x) f'(a) - \frac{B_2}{4!} (\delta x)^2 f''(a) + \dots \right]$$

or

$$\sum_{x=a}^{x=a+(n-1)\delta x} f(x) = \frac{1}{\delta x} \int_a^{a+n\delta x} f(x) dx - \frac{1}{2} [f(a + n\delta x) - f(a)]$$

$$+ \frac{B_1}{2!} (\delta x) [f'(a + n\delta x) - f'(a)] - \frac{B_2}{4!} (\delta x)^2 [f''(a + n\delta x) - f''(a)] + \dots$$

Finally, on rearrangement, we obtain:

$$\int_a^{a+n\delta x} f(x) dx = \frac{(\delta x)}{2} [f(a) + 2f(a + \delta x) + 2f(a + 2\delta x) + \dots$$

$$+ 2f(a + (n-1)\delta x) + f(a + n\delta x)]$$

$$- \frac{(\delta x)^2}{12} [f'(a + n\delta x) - f'(a)] + \frac{(\delta x)^4}{720} [f'''(a + n\delta x) - f'''(a)]$$

$$- \frac{(\delta x)^6}{30240} [f^{(5)}(a + n\delta x) - f^{(5)}(a)] + \dots (4.4.5)$$

which is the normal form of the expansion.

It may be mentioned that the Euler-Maclaurin formula is of little practical use as an integration procedure. It finds an application, however, in the summation of series and, in this connection, will receive attention in Chapter 5.

4.5 CENTRAL DIFFERENCE FORMULAE FOR INTEGRATION

The Bickley expansions (equations 4.2.8 and 4.2.9) can be used to obtain several useful formulae for integration which involve central differences. Thus, from equation 4.2.8 with $n = -2$ we obtain:

$$\left(\frac{\delta}{U} \right)^2 = 1 + \frac{1}{12} \delta^2 - \frac{1}{240} \delta^4 + \frac{31}{60480} \delta^6 - \dots (4.5.1)$$

which may be written:

$$\delta^2 \cdot \frac{1}{D^2}(f_0) = (\delta x)^2 (f_0 + \frac{1}{12}\delta^2 f_0 - \frac{1}{240}\delta^4 f_0 + \frac{31}{80480}\delta^6 f_0 - \dots) \quad (4.5.2)$$

Now putting $f_0 = g''_0 = D^2(g_0)$ equation 4.5.2 becomes:

$$\delta^2 g_0 = (\delta x)^2 (g''_0 + \frac{1}{12}\delta^2 g''_0 - \frac{1}{240}\delta^4 g''_0 + \frac{31}{80480}\delta^6 g''_0 - \dots) \quad (4.5.3)$$

which expresses the second central difference in terms of the central difference of the second derivative at the point in question, and has several applications to the solution of differential equations.

Again, if we put $n = -1$ in equation 4.2.9, there results:

$$\frac{1}{\mu} \left(\frac{\delta}{U} \right) = 1 - \frac{1}{12}\delta^2 + \frac{1}{720}\delta^4 - \frac{1}{60480}\delta^6 + \dots \quad (4.5.4)$$

or

$$\delta \cdot \frac{1}{D}(f_0) = (\delta x) [\mu f_0 - \frac{1}{12}\delta^2(\mu f_0) + \frac{1}{720}\delta^4(\mu f_0) - \frac{1}{60480}\delta^6(\mu f_0) + \dots] \quad (4.5.5)$$

Now, replace f_0 by $g'_1 = D(g_1)$, and observe that:

$$\left. \begin{aligned} \delta g_1 &= \delta E^{\frac{1}{2}} g_0 = (E^{\frac{1}{2}} - E^{-\frac{1}{2}}) E^{\frac{1}{2}} g_0 = (E - 1) g_0 = g_1 - g_0 \\ \mu g_1 &\equiv \mu E^{\frac{1}{2}} g_0 \equiv \frac{1}{2}(E^{\frac{1}{2}} + E^{-\frac{1}{2}}) E^{\frac{1}{2}} g_0 \equiv \frac{1}{2}(E + 1) g_0 \equiv \frac{1}{2}(g_1 + g_0) \end{aligned} \right\} \dots \quad (4.5.6)$$

whence, equation 4.5.5 becomes:

$$g_1 - g_0 = \frac{1}{2}(\delta x) [g'_1 + g'_0 - \frac{1}{12}(\delta^2 g'_1 + \delta^2 g'_0) + \frac{1}{720}(\delta^4 g'_1 + \delta^4 g'_0) - \frac{1}{60480}(\delta^6 g'_1 + \delta^6 g'_0) - \dots] \quad (4.5.7)$$

This is clearly an integration formula expressing

$$g_1 - g_0 \left(= \int_0^1 g'(x) dx \right)$$

in terms of the relevant functional values of $g'(x)$.

Finally we shall consider the expansion of the operator $(\mu\delta)$.

We have:

$$(\mu\delta)f_0 = \frac{1}{2}(E^{\frac{1}{2}} + E^{-\frac{1}{2}})(E^{\frac{1}{2}} - E^{-\frac{1}{2}})f_0 = \frac{1}{2}(E^1 - E^{-1})f_0 = \frac{1}{2}(f_1 - f_{-1}) \quad (4.5.8)$$

Again, equation 4.5.4 may be written:

$$\mu \left(\frac{\delta}{U} \right) = \mu^2 (1 - \frac{1}{12}\delta^2 + \frac{1}{720}\delta^4 - \frac{1}{60480}\delta^6 + \dots)$$

or, observing that $\mu^2 = 1 + \frac{1}{4}\delta^2$ (equation 4.2.5)

$$\frac{\mu\delta}{U} = 1 + \frac{1}{8}\delta^2 - \frac{1}{180}\delta^4 + \frac{1}{1512}\delta^6 - \dots \quad (4.5.9)$$

It follows that:

$$\frac{1}{2} \frac{(f_1 - f_{-1})}{U} = f_0 + \frac{1}{8}\delta^2 f_0 - \frac{1}{180}\delta^4 f_0 + \frac{1}{1512}\delta^6 f_0 - \dots$$

or, writing: $f_x = g'_x = Dg_x$

$$\mu\delta g_0 = \frac{1}{2}(g_1 - g_{-1}) = (\delta x) (g'_0 + \frac{1}{8}\delta^2 g'_0 - \frac{1}{180}\delta^4 g'_0 + \frac{1}{1512}\delta^6 g'_0 - \dots) \quad (4.5.10)$$

which is the required result.

4.6 THE CHEBYSHEV INTEGRATION FORMULAE

Integration formulae of the Simpson-Weddle type are advantageous when used in conjunction with a table of functional values, since they make use of equal intervals. On the other hand, each of the function values has, in general, to be multiplied by a different coefficient. When function values have to be calculated for the purpose of the integration, it is often more convenient to use an integration formula in which all such values are treated in the same manner.

We thus examine the possibility of finding an integration formula:

$$\int_a^b g(x) dx = k[g(x_1) + g(x_2) + g(x_3) + \dots + g(x_n)] + R_n \quad (4.6.1)$$

in which k and x_1, x_2, \dots, x_n are independent of the particular function $g(x)$ being integrated.

Actually, a slightly more general result, given by MILNE-THOMSON⁽³⁾ will be examined:

$$\int_{-1}^{+1} F(x)g(x) dx = k[g(x_1) + g(x_2) + \dots + g(x_n)] + R_n \quad (4.6.2)$$

where k, x_1, x_2, \dots, x_n depend on $F(x)$ but not on $g(x)$.

First, assume that $g(x)$ can be expanded in a Maclaurin series:

$$g(x) = g(0) + xg'(0) + \frac{x^2}{2!}g''(0) + \dots + \frac{x^n}{n!}g^{(n)}(0) + \frac{x^{n+1}}{(n+1)!}g^{(n+1)}(\xi) \quad (4.6.3)$$

in which the remainder term has $0 \leq \xi \leq x$.

We then have:

$$\begin{aligned} \int_{-1}^{+1} F(x)g(x)dx &= \int_{-1}^{+1} F(x) \left[g(0) + xg'(0) + \frac{x^2}{2!}g''(0) + \dots \right. \\ &\quad \left. + \frac{x^n}{n!}g^{(n)}(0) + \frac{x^{n+1}}{(n+1)!}g^{(n+1)}(\xi) \right] dx \\ &= \sum_{m=0}^n \frac{g^{(m)}(0)}{m!} \int_{-1}^{+1} F(x) \cdot x^m dx + \int_{-1}^{+1} \frac{g^{(n+1)}(\xi)}{(n+1)!} x^{n+1} F(x) dx \quad \dots (4.6.4) \end{aligned}$$

Again, using equation 4.6.3 with $x = x_1, x_2, \dots, x_n$:

$$\begin{aligned} k[g(x_1) + g(x_2) + \dots + g(x_n)] &= kng(0) + kg'(0) \sum_{r=1}^n x_r \\ &\quad + \frac{kg''(0)}{2!} \sum_{r=1}^n x_r^2 + \dots + \frac{kg^{(n)}(0)}{n!} \sum_{r=1}^n x_r^n + \frac{k}{(n+1)!} \sum_{r=1}^n x_r^{n+1} g^{(n+1)}(\xi_r) \quad \dots (4.6.5) \end{aligned}$$

Whence, if equations 4.6.4 and 4.6.5 are to be identical to the error specified by R_n in equation 4.6.2, we have, on equating coefficients of $g(0), g'(0), \dots, g^{(n)}(0)$:

$$\left. \begin{aligned} kn &= \int_{-1}^{+1} F(x) dx \\ k \sum_{r=1}^n x_r &= \int_{-1}^{+1} xF(x) dx \\ k \sum_{r=1}^n x_r^2 &= \int_{-1}^{+1} x^2 F(x) dx \\ &\dots \\ k \sum_{r=1}^n x_r^n &= \int_{-1}^{+1} x^n F(x) dx \end{aligned} \right\} \quad \dots (4.6.6)$$

Together with:

$$R_n = \int_{-1}^{+1} \frac{g^{(n+1)}(\xi)}{(n+1)!} x^{n+1} F(x) dx - \frac{k}{(n+1)!} \sum_{r=1}^n x_r^{n+1} g^{(n+1)}(\xi_r) \quad \dots (4.6.7)$$

where $0 \leq \xi \leq x$ and $0 \leq \xi_r \leq x_r$ ($r = 1, \dots, n$).

To determine the values of x_1, x_2, \dots, x_n from equation 4.6.6 we may adopt the following procedure.

First let:

$$\sum_{r=1}^n x_r^m = s_m \quad \dots (4.6.8)$$

$$\text{so that} \quad s_m = \frac{1}{k} \int_{-1}^{+1} x^m F(x) dx = \frac{n \int_{-1}^{+1} x^m F(x) dx}{\int_{-1}^{+1} F(x) dx} \quad \dots (4.6.9)$$

The quantities s_m are thus the sums of the m th powers of the n quantities x_1, x_2, \dots, x_n which suggests the roots of an equation of degree n . Assume then that x_1, x_2, \dots, x_n are the roots of:

$$f_n(x) \equiv x^n + a_1 x^{n-1} + a_2 x^{n-2} + \dots + a_n = 0 \quad \dots (4.6.10)$$

$$\begin{aligned} \text{so that:} \quad x^n + a_1 x^{n-1} + a_2 x^{n-2} + \dots + a_n \\ \equiv (x - x_1)(x - x_2) \dots (x - x_n) \quad \dots (4.6.11) \end{aligned}$$

Now put $x = 1/y$ so that, multiplying both sides by y^n equation 4.6.11 becomes:

$$1 + a_1 y + a_2 y^2 + \dots + a_n y^n \equiv (1 - x_1 y)(1 - x_2 y) \dots (1 - x_n y)$$

Taking logs of both sides:

$$\begin{aligned} \log_e(1 + a_1 y + a_2 y^2 + \dots + a_n y^n) &= \sum_{r=1}^n \log_e(1 - x_r y) \\ &= - \sum_{r=1}^n (x_r y + \frac{x_r^2 y^2}{2} + \frac{x_r^3 y^3}{3} + \dots + \frac{x_r^m y^m}{m} + \dots) \\ &= -s_1 y - \frac{s_2}{2} y^2 - \frac{s_3}{3} y^3 - \dots - s_m y^m - \dots \quad \dots (4.6.12) \end{aligned}$$

on expansion and by virtue of equation 4.6.8.

Next take exponentials of both sides:

$$\begin{aligned} 1 + a_1 y + a_2 y^2 + \dots + a_n y^n \\ = \exp[-s_1 y - (s_2/2)y^2 - (s_3/3)y^3 - \dots - (s_m/m)y^m - \dots] \quad \dots (4.6.13) \end{aligned}$$

So that, upon expanding the exponential, and equating coefficients, we can determine the coefficients a_m in terms of the power sums s_m . The first few are:

$$\begin{aligned} a_1 &= -s_1 \\ a_2 &= -\frac{1}{2}s_2 + \frac{1}{2}s_1^2 \\ a_3 &= -\frac{1}{3}s_3 + \frac{1}{2}s_1 s_2 - \frac{1}{6}s_1^3 \\ a_4 &= -\frac{1}{4}s_4 + \frac{1}{2}s_1 s_3 - \frac{1}{4}s_1^2 s_2 + \frac{1}{8}s_2^2 + \frac{1}{24}s_1^4 \\ a_5 &= -\frac{1}{5}s_5 + \frac{1}{4}s_1 s_4 + \frac{1}{6}s_2 s_3 - \frac{1}{6}s_1^2 s_3 - \frac{1}{8}s_1 s_2^2 + \frac{1}{12}s_1^3 s_2 - \frac{1}{120}s_1^5 \\ a_6 &= -\frac{1}{6}s_6 + \frac{1}{5}s_1 s_5 + \frac{1}{6}s_2 s_4 + \frac{1}{18}s_3^2 - \frac{1}{8}s_1^2 s_4 - \frac{1}{6}s_1 s_3 s_2 \\ &\quad - \frac{1}{48}s_2^3 + \frac{1}{18}s_1^2 s_3 + \frac{1}{16}s_1^2 s_2^2 - \frac{1}{48}s_1^4 s_2 + \frac{1}{720}s_1^6 \end{aligned}$$

Thus, knowing the form of $F(x)$ and using equation 4.6.9 to determine the values of s_m , we can write down the equation 4.6.10 and, by solving it, obtain the values of $x_1 \dots x_n$.

The most usual technique is to take $F(x) = 1$, we thus obtain the following equations $f_n(x) = 0$, and the relevant roots for insertion in 4.6.2.

n	$k(=2/n)$	$f_n(x)$	Roots
2	1	$x^2 - \frac{1}{2}$	$\pm .5773\ 503$
3	$\frac{2}{3}$	$x(x^2 - \frac{1}{2})$	$0, \pm .7071\ 068$
4	$\frac{1}{2}$	$x^4 - \frac{3}{8}x^2 + \frac{1}{8}$	$\pm .1875\ 925 \pm .7946\ 545$
5	$\frac{2}{5}$	$x(x^4 - \frac{5}{8}x^2 + \frac{1}{8})$	$0, \pm .3745\ 414 \pm .8324\ 975$
6	$\frac{1}{3}$	$x^6 - x^4 + \frac{1}{3}x^2 - \frac{1}{108}$	$\pm .2666\ 353 \pm .4225\ 186, \pm .8662\ 476$

It should be noticed that the remainder, after an n point integration using the Chebyshev method, is zero if $g(x)$ is a polynomial of degree not greater than n .

The Chebyshev method is particularly suited to use on an automatic digital computer since, in such machines, the fewer different kinds of arithmetic operations which have to be performed to arrive at a given result the better.

[An alternative method of approach, which gives the polynomials $f_n(x)$ directly, is as follows:

Let

$$f_n(z) \equiv (z - x_1)(z - x_2) \dots (z - x_n) = 0$$

then:

$$\begin{aligned} f_n(z) &= z^n \left(1 - \frac{x_1}{z}\right) \left(1 - \frac{x_2}{z}\right) \dots \left(1 - \frac{x_n}{z}\right) \\ &= z^n \exp \sum_{r=1}^n \log \left(1 - \frac{x_r}{z}\right) \\ &= z^n \exp \sum_{r=1}^n \left(-\sum_{t=1}^{\infty} \frac{x_r^t}{t z^t}\right) \end{aligned}$$

Since n is finite, we may interchange the summation symbols so that:

$$f_n(z) = z^n \exp - \sum_{t=1}^{\infty} \left(\frac{\sum_{r=1}^n x_r^t}{t z^t}\right)$$

but, by equation 4.6.9,

$$\sum_{r=1}^n x_r^t = n \int_{-1}^{+1} x^t F(x) dx / \int_{-1}^{+1} F(x) dx$$

whence,

$$\begin{aligned} f_n(z) &= z^n \exp - \sum_{t=1}^{\infty} n \int_{-1}^{+1} x^t F(x) dx / t z^t \int_{-1}^{+1} F(x) dx \\ &= z^n \exp n \int_{-1}^{+1} F(x) \left[-\sum_{t=1}^{\infty} \frac{x^t}{t z^t} \right] dx / \int_{-1}^{+1} F(x) dx \\ &= z^n \exp n \int_{-1}^{+1} F(x) \log \left(1 - \frac{x}{z}\right) dx / \int_{-1}^{+1} F(x) dx \end{aligned}$$

where $f(z)$ is taken to include only positive powers of z (including zero) in the expansion. When $F(x) = 1$ this becomes very simply:

$$f_n(z) = z^n \exp - n \left(\frac{1}{2.3z^2} + \frac{1}{4.5z^4} + \frac{1}{6.7z^6} + \dots \right)$$

from which the preceding values of $f_n(x)$ can be quickly derived by expansion.]

It may be noted, in conclusion, that the Chebyshev integration formula becomes unsatisfactory for $n = 8$ and for $n > 9$ since, at this point, the values of x_r are no longer restricted to the range $(-1, +1)$ and, in fact, may not even be real.

4.7 THE GAUSSIAN INTEGRATION FORMULA

The Chebyshev method of approximate integration is one of a family of such results which may be written:

$$\int_a^b F(x) \cdot g(x) dx = w_1 \cdot g(x_1) + w_2 \cdot g(x_2) \dots + w_n g(x_n) + R_n \dots (4.7.1)$$

where the weight functions* (w_r) and the associated points (x_r) are independent of the particular function $g(x)$ which is the subject of the integration. [They depend, however, on $F(x)$.]

In section 4.6 we investigated the way in which the values (x_r) had to be chosen in order to make the weight functions equal. We shall now consider the problem in slightly more general terms. First take $F(x) = 1$ and assume that $g(x)$ can be expanded in a convergent power series in x so that:

$$g(x) = a_0 + a_1 x + \dots + a_n x^n + \dots \dots (4.7.2)$$

* Sometimes known as Christoffel numbers.

Next, take any set of n points $x_1 \dots x_n$ and set up the Lagrangean approximation to $g(x)$ (see section 3.7). This may be written:

$$L(x) = \prod_n(x) \sum_{r=1}^n \frac{g(x_r)}{\prod_n'(x_r) \cdot (x - x_r)} \quad \dots (4.7.3)$$

where

$$\prod_n(x) = (x - x_1)(x - x_2) \dots (x - x_n) \quad \dots (4.7.4)$$

and $\prod_n'(x)$ is used to represent $d\prod_n/dx$.

Now $L(x)$ coincides in value with the function $g(x)$ at each of the points (x_r) ($r = 1 \dots n$) so that we may write:

$$g(x) - L(x) = \prod_n(x) (b_0 + b_1x + \dots + b_sx^s + \dots)$$

where the series on the right is again assumed to be convergent. We thus obtain:

$$\begin{aligned} \int_a^b g(x) dx &= \int_a^b L(x) dx + \int_a^b \left(\prod_n(x) \sum_{s=0}^{\infty} b_s x^s \right) dx \\ &= \sum_{r=1}^n w_r g(x_r) + R_n \end{aligned} \quad \dots (4.7.5)$$

where

$$w_r = \int_a^b \frac{\prod_n(x) dx}{(x - x_r) \prod_n'(x_r)} \quad \dots (4.7.6)$$

and

$$R_n = \int_a^b \left(\prod_n(x) \sum_{s=0}^{\infty} b_s x^s \right) dx \quad \dots (4.7.7)$$

Thus, if $g(x)$ is a polynomial, of maximum degree n , we see that $R_n = 0$, and that the n values (x_r) can be chosen arbitrarily so long as we choose w_r as defined by equation 4.7.6.

Since this situation is equivalent to saying that there are arbitrary constants available, it is reasonable to enquire if these can be determined so as to make equation 4.7.5, with $R_n = 0$, true when $g(x)$ is any polynomial of degree $(2n - 1)$.

Thus we require

$$R_n = \int_a^b \left[\prod_n(x) \sum_{s=0}^{n-1} b_s x^s \right] dx = 0$$

for arbitrary b_s ($s \leq n - 1$), and this can only be true if:

$$\int_a^b \prod_n(x) \cdot x^s dx = 0 \quad (s = 0, 1, \dots, n - 1) \quad \dots (4.7.8)$$

Now it is well known that the Legendre polynomials, which are conveniently defined for this purpose by means of Rodrigues formula⁽⁴⁾

$$P_0(x) = 1, P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n \quad \dots (4.7.9)$$

are such that⁽⁵⁾:

$$\int_{-1}^{+1} P_n(x) \cdot x^s dx = 0$$

for all integer values of s less than n ; furthermore they have n distinct real roots in the interval $(-1, +1)$.

The integral (4.7.8) can be transformed by means of the substitution:

$$x = \frac{b-a}{2} y + \frac{b+a}{2}$$

to take the limits $(-1, +1)$, and it is thus seen that if, after the transformation, we take the values of (y_r) to be the roots of $P_n(y) = 0$ the resulting n point quadrature formula will be exact for all polynomials $g(x)$ whose degree is less than $(2n)$.

An identical argument can be carried out for the more general integral (4.7.1). If this is assumed to be transformed to limits $(-1, +1)$ we have:

$$\int_{-1}^{+1} F(x) g(x) dx = w_1 \cdot g(x_1) + w_2 \cdot g(x_2) + \dots + w_n g(x_n) + R_n \quad \dots (4.7.10)$$

where:

$$w_r = \frac{F(x_r)}{\prod_n'(x_r)} \int_{-1}^{+1} \frac{\prod_n(x)}{(x - x_r)} dx \quad \dots (4.7.11)$$

and

$$R_n = \int_{-1}^{+1} \left(\prod_n(x) \sum_{s=0}^{\infty} b_s x^s \right) dx \quad \dots (4.7.12)$$

(x_r) ($r = 1 \dots n$) being the roots of $P_n(x) = 0$.

The values of the roots (x_r) , and of the corresponding weight factors w_r , for the important case $F(x) = 1$ are given in Table

4.7.1. They are taken from the basic table of LOWAN, DAVIDS and LEVENSON⁽⁶⁾ who give values to 15 decimal places for values of n up to 16.

Table 4.7.1. Gaussian integration points and coefficients

$x_1 = 0$	$n = 1$	$w_1 = 2$
$x_1 = -x_2 = .57735\ 02692$	$n = 2$	$w_1 = w_2 = 1$
$x_1 = -x_3 = .77459\ 66692$ $x_2 = 0$	$n = 3$	$w_1 = w_3 = \frac{8}{9}$ $w_2 = \frac{8}{9}$
$x_1 = -x_4 = .86113\ 63116$ $x_2 = -x_3 = .33998\ 10436$	$n = 4$	$w_1 = w_4 = .34785\ 48451$ $w_2 = w_3 = .65214\ 51549$
$x_1 = -x_5 = .90617\ 98459$ $x_2 = -x_4 = .53846\ 93101$ $x_3 = 0$	$n = 5$	$w_1 = w_5 = .23692\ 68851$ $w_2 = w_4 = .47862\ 86705$ $w_3 = .56888\ 88889$
$x_1 = -x_6 = .93246\ 95142$ $x_2 = -x_5 = .66120\ 93865$ $x_3 = -x_4 = .23861\ 91861$	$n = 6$	$w_1 = w_6 = .17132\ 44924$ $w_2 = w_5 = .36076\ 15730$ $w_3 = w_4 = .46791\ 39346$
$x_1 = -x_7 = .94910\ 79123$ $x_2 = -x_6 = .74153\ 11856$ $x_3 = -x_5 = .40584\ 51514$ $x_4 = 0$	$n = 7$	$w_1 = w_7 = .12948\ 49662$ $w_2 = w_6 = .27970\ 53915$ $w_3 = w_5 = .38183\ 00505$ $w_4 = .41795\ 91837$
$x_1 = -x_8 = .96028\ 98565$ $x_2 = -x_7 = .79666\ 64774$ $x_3 = -x_6 = .52553\ 24099$ $x_4 = -x_5 = .18343\ 46425$	$n = 8$	$w_1 = w_8 = .10122\ 85363$ $w_2 = w_7 = .22238\ 10345$ $w_3 = w_6 = .31370\ 66459$ $w_4 = w_5 = .36268\ 37834$
$x_1 = -x_9 = .96816\ 02395$ $x_2 = -x_8 = .83603\ 11073$ $x_3 = -x_7 = .61337\ 14327$ $x_4 = -x_6 = .32425\ 34234$ $x_5 = 0$	$n = 9$	$w_1 = w_9 = .08127\ 43884$ $w_2 = w_8 = .18064\ 81607$ $w_3 = w_7 = .26061\ 06964$ $w_4 = w_6 = .31234\ 70770$ $w_5 = .33023\ 93550$
$x_1 = -x_{10} = .97390\ 65285$ $x_2 = -x_9 = .86506\ 33667$ $x_3 = -x_8 = .67940\ 95683$ $x_4 = -x_7 = .43339\ 53941$ $x_5 = -x_6 = .14887\ 43399$	$n = 10$	$w_1 = w_{10} = .06667\ 13443$ $w_2 = w_9 = .14985\ 13492$ $w_3 = w_8 = .21908\ 63625$ $w_4 = w_7 = .26926\ 67193$ $w_5 = w_6 = .29552\ 42247$

It may be noted that, $\sum_{r=1}^n w_r = 2$ —a result which is obtained by taking $F(x) = g(x) = 1$ in equation 4.7.10.

The Gaussian integration method just described may be extended to other ranges than $(-1, +1)$. Probably the most important from the practical point of view are those ranges in which either, or both of (a, b) are infinite. Finiteness of the resulting integral implies that, if $g(x)$ is a polynomial, $F(x)$ is not, and it is natural

to relate the extended method to particular forms of $F(x)$. Thus consider the integral:

$$\int_0^\infty e^{-x} g(x) dx \quad \dots (4.7.13)$$

where $g(x)$ is again expressed by equation 4.7.2. We once more seek an integration formula of the type shown in equation 4.7.1 and set up a Lagrangean approximation, but this time to $g(x)$ only. If:

$$g(x) - L(x) = \prod_n (x - x_r) (b_0 + b_1 x + \dots + b_s x^s + \dots)$$

we obtain, from equation 4.7.13

$$\begin{aligned} \int_0^\infty e^{-x} g(x) dx &= \int_0^\infty e^{-x} L(x) dx + \int_0^\infty \left[\prod_n (x - x_r) \sum_{s=0}^\infty b_s x^s \right] e^{-x} dx \\ &= \sum_{r=1}^n w_r g(x_r) + R_n \quad \dots (4.7.14) \end{aligned}$$

where

$$w_r = \frac{1}{\prod_n (x_r - x_s)} \int_0^\infty \frac{e^{-x} \prod_n (x - x_s)}{(x - x_r)} dx \quad \dots (4.7.15)$$

and

$$R_n = \int_0^\infty \left[\prod_n (x - x_r) \sum_{s=0}^\infty b_s x^s \right] e^{-x} dx \quad \dots (4.7.16)$$

If we wish equation 4.7.14 to be exact for all polynomials $g(x)$ of degree less than $(2n - 1)$ we thus require:

$$\int_0^\infty \prod_n (x - x_r) x^s e^{-x} dx = 0 \quad (s = 0, 1, \dots, n - 1)$$

Now in this case the Laguerre polynomials defined by:

$$L_n(x) = e^x \frac{d^n}{dx^n} (x^n e^{-x}) \quad \dots (4.7.17)$$

are appropriate since it can be shown⁽⁷⁾ that:

$$\int_0^\infty e^{-x} x^s L_n(x) dx = 0$$

whenever $s < n$, and we thus take the sampling points (x_r) to be the roots of $L_n(x) = 0$. Some weights and coefficients are given⁽⁹⁾ in Table 4.7.2.

Table 4.7.2. Laguerre polynomial roots and weight coefficients

$x_1 = 1$	$n = 1$	$w_1 = 1$
$x_1 = 0.58578\ 64$ $x_2 = 3.41421\ 36$	$n = 2$	$w_1 = .85355\ 34$ $w_2 = .14644\ 66$
$x_1 = 0.41577\ 46$ $x_2 = 2.29428\ 04$ $x_3 = 6.28994\ 51$	$n = 3$	$w_1 = .71109\ 30$ $w_2 = .27851\ 77$ $w_3 = .01038\ 93$
$x_1 = 0.32254\ 77$ $x_2 = 1.74576\ 11$ $x_3 = 4.53662\ 03$ $x_4 = 9.39507\ 09$	$n = 4$	$w_1 = .60315\ 41$ $w_2 = .35741\ 87$ $w_3 = .03888\ 79$ $w_4 = .00053\ 93$
$x_1 = 0.26356\ 03$ $x_2 = 1.41340\ 31$ $x_3 = 3.59642\ 58$ $x_4 = 7.08581\ 00$ $x_5 = 12.64080\ 08$	$n = 5$	$w_1 = .52175\ 56$ $w_2 = .39866\ 68$ $w_3 = .07594\ 24$ $w_4 = .00361\ 18$ $w_5 = .00002\ 34$

Again, if the integral required is of the form:

$$\int_{-\infty}^{+\infty} e^{-x^2} g(x) dx$$

it is appropriate to consider the Hermite polynomials⁽⁸⁾ defined by:

$$H_n(x) = (-1)^n e^{x^2} \frac{d^n}{dx^n} (e^{-x^2}) \quad (4.7.18)$$

for which it can be shown that:

$$\int_{-\infty}^{+\infty} e^{-x^2} x^s H_n(x) dx = 0$$

whenever $s < n$. These lead to the formula

$$\int_{-\infty}^{+\infty} e^{-x^2} g(x) dx = \sum_{r=1}^n w_r g(x_r) + R_n \quad \dots (4.7.19)$$

where

$$w_r = \frac{1}{H'_n(x_r)} \int_{-\infty}^{+\infty} \frac{e^{-x^2} H_n(x)}{(x - x_r)} dx$$

$$R_n = \int_{-\infty}^{+\infty} \left[\frac{H_n(x)}{n} \sum_{s=0}^{\infty} b_s x^s \right] e^{-x^2} dx$$

and (x_r) are the roots of $H_n(x) = 0$.

In a like manner the Chebyshev polynomials:

$$T_n(x) = \frac{1}{2^{n-1}} \cos(n \cos^{-1} x) \quad \dots (4.7.20)$$

have the property⁽¹⁰⁾:

$$\int_{-1}^{+1} \frac{x^s T_n(x)}{\sqrt{1-x^2}} dx = 0$$

for $s < n$ which leads at once to:

$$\int_{-1}^{+1} \frac{g(x) dx}{\sqrt{1-x^2}} = \sum_{r=1}^n w_r g(x_r) + R_n \quad \dots (4.7.21)$$

where

$$w_r = \frac{1}{H'_n(x_r)} \int_{-1}^{+1} \frac{\Pi(x)}{(x - x_r) \sqrt{1-x^2}} dx$$

$$R_n = \int_{-1}^{+1} \left[\frac{\Pi(x)}{n} \sum_{s=0}^{\infty} b_s x^s \right] \frac{dx}{\sqrt{1-x^2}}$$

and the (x_r) are the roots of $T_n(x) = 0$.

 Table 4.7.3. Roots and weight coefficients for $(1-x^2) P'_n(x)$

$x_1 = -x_2 = 1$	$n = 1$	$w_1 = w_2 = 1$
$x_1 = -x_2 = 1$ $x_2 = 0$	$n = 2$	$w_1 = w_3 = 1/3$ $w_2 = 4/3$
$x_1 = -x_4 = 1$ $x_3 = -x_3 = 0.44721\ 36$	$n = 3$	$w_1 = w_4 = 0.16666\ 67$ $w_2 = w_3 = 0.83333\ 33$
$x_1 = -x_5 = 1$ $x_2 = -x_4 = 0.65465\ 37$ $x_3 = 0$	$n = 4$	$w_1 = w_5 = 0.10000\ 00$ $w_2 = w_4 = 0.54444\ 44$ $w_3 = 0.71111\ 11$
$x_1 = -x_6 = 1$ $x_2 = -x_5 = 0.76505\ 53$ $x_3 = -x_4 = 0.08135\ 70$	$n = 5$	$w_1 = w_6 = 0.66666\ 67$ $w_2 = w_5 = 0.37847\ 50$ $w_3 = w_4 = 0.55485\ 84$
$x_1 = -x_7 = 1$ $x_2 = -x_6 = 0.83022\ 39$ $x_3 = -x_5 = 0.46884\ 88$ $x_4 = 0$	$n = 6$	$w_1 = w_7 = 0.04761\ 90$ $w_2 = w_6 = 0.27682\ 60$ $w_3 = w_5 = 0.43174\ 54$ $w_4 = 0.48761\ 90$

To conclude this section it may be mentioned that it is sometimes desirable to have an integration formula of Gaussian type, but involving the function values at the limits of the range of integration (a, b). Such results are usually known as 'Lobatto's or Radau's formulae' and a typical example results from the use of the polynomial:

$$(1 - x^2)P'_n(x) \quad \dots (4.7.22)$$

instead of equation 4.7.9. Formulae based on equation 4.7.22 are less accurate than the true Gaussian results, being exact only for polynomials of degree $(2n - 1)$ when $(n + 1)$ points are used. The chief application is to the integration of functions having the value zero at both ends of the range of integration.

Roots and multiples for the polynomials 4.7.22 are given⁽¹¹⁾ in Table 4.7.3.

4.8 OSCILLATING INTEGRANDS

A class of integrals which is of practical importance is that typified by:

$$\int_a^b f(p) \frac{\sin}{\cos} (xp) dp$$

For finite values of a and b Filon's method⁽¹²⁾ is very suitable. Here, the range is divided into $2n$ equal parts so that

$$b = a + 2nh$$

and, if

$$\theta = hx$$

it can be shown that:

$$\int_a^b f(p) \sin (xp) dp = h[- a\{f(b) \cos xb - f(a) \cos xa\} + \beta S_{2n} + \gamma S_{2n-1}]$$

$$\int_a^b f(p) \cos (xp) dp = h[a\{f(b) \sin xb - f(a) \sin xa\} + \beta C_{2n} + \gamma C_{2n-1}]$$

α, β and γ are defined by:

$$\alpha = \frac{1}{\theta^3} [\theta^2 + \theta \sin \theta \cos \theta - 2 \sin^2 \theta]$$

$$\beta = \frac{2}{\theta^3} [\theta(1 + \cos^2 \theta) - 2 \sin \theta \cos \theta]$$

$$\gamma = \frac{4}{\theta^3} [\sin \theta - \theta \cos \theta]$$

S_{2n} and C_{2n} are the sums of all even ordinates of the curves $y = f(p) \frac{\sin}{\cos} (xp)$ between a and b inclusive less one-half of the first and last ordinates. S_{2n-1} and C_{2n-1} are the sums of all odd ordinates. Tables of α, β and γ have been produced by Filon.

For an infinite upper limit, the method of choice is to use:

$$\int_{n\pi}^{\infty} f(x) \sin x dx = (-1)^n (1 + a_1 \delta^2 + a_2 \delta^4 + \dots) f(n\pi)$$

where the a_i are the coefficients $(\delta^2)^i$ in the expansion of $(1 + (\log E)^2/\pi^2)^{-1}$ and:

$$a_1 = -0.10132 \ 118$$

$$a_2 = +0.01870 \ 941$$

$$a_3 = -0.00387 \ 695$$

$$a_4 = +0.00084 \ 579$$

etc.

4.9 ERRORS OF FINITE DIFFERENCE APPROXIMATION

Except in the case of the Chebyshev formula, we have not attempted to estimate the errors produced by the approximate integration formulae. A general method of arriving at such error estimates has been given by MILNE⁽¹³⁾, and we shall now indicate its application by considering the error produced by Simpson's rule (4.3.4) and by its Newton-Cotes *open* variant (4.3.10). These particular examples have been chosen because the error estimates will be needed later, in Chapter 6.

In the first place, we make the observation that all finite difference formulae are exact when applied to any polynomial whose degree, n , does not exceed a certain upper bound, N , say, which is a constant for the particular formula involved. It follows that if the error is to be estimated when the finite difference formula is applied to an arbitrary function, a suitable technique will be to expand that function in a Maclaurin series:

$$f(x) = f(0) + \frac{x}{1!} f'(0) + \frac{x^2}{2!} f''(0) + \dots + \frac{x^n}{n!} f^{(n)}(0) + \frac{x^{n+1}}{(n+1)!} f^{(n+1)}(\xi) \quad \dots (4.9.1)$$

where $0 \leq \xi \leq x$.

Suppose that the *actual* operation to be performed is O and that the finite difference process is represented by F ; the error, or remainder R , is then given by:

$$R = O[f(x)] - F[f(x)] \quad \dots (4.9.2)$$

Now, if in equation 4.9.1 we make $n = N$, where N is the highest degree polynomial for which F is exact, we may write 4.9.2 as:

$$R_N = O\left[\frac{x^{N+1}}{(N+1)!} f^{(N+1)}(\xi)\right] - F\left[\frac{x^{N+1}}{(N+1)!} f^{(N+1)}(\xi)\right] \quad \dots (4.9.3)$$

(Since O and F are assumed to be *linear* operators, so that

$$O(f_1 + f_2 + \dots + f_n) = O(f_1) + O(f_2) + \dots + O(f_n) \text{ etc.})$$

whence, replacing $f^{(N+1)}(\xi)$ by its maximum value in the interval:

$$R_N \leq \frac{1}{(N+1)!} f_{\max}^{(N+1)} \{O(x^{N+1}) - F(x^{N+1})\} \quad \dots (4.9.4)$$

We now apply this result to estimate the error produced by Simpson's rule. We know that the result is exact for quadratic polynomials, let us therefore try $N = 2$ in equation 4.9.4.

$$O = \int_0^{2(\delta x)} x^3 dx = 4(\delta x)^4$$

$$F = \frac{(\delta x)}{3} [f(2) + 4f(1) + f(0)] = \frac{(\delta x)}{3} [(2\delta x)^3 + 4(\delta x)^3] = 4(\delta x)^4$$

whence the rule is also exact for cubic polynomials.

Next put $N = 3$

$$O = \int_0^{2\delta x} x^4 dx = 32 \frac{(\delta x)^5}{5}$$

$$F = \frac{(\delta x)}{3} [(2\delta x)^4 + 4(\delta x)^4] = 20 \frac{(\delta x)^5}{3}$$

Whence, from equation 4.9.4:

$$R_N \leq \frac{1}{4!} f_{\max}^{(4)} \left(\frac{32}{5} - \frac{20}{3}\right) (\delta x)^5 = -\frac{(\delta x)^5}{90} f_{\max}^{(4)} \quad \dots (4.9.5)$$

which is the required estimate.

Next consider the Newton-Cotes formula (4.3.10).

Here again it is readily shown that for cubic polynomials the formula is exact. For $N = 3$ we obtain:

$$O = \int_0^{4(\delta x)} x^4 dx = 1024 \frac{(\delta x)^5}{5}$$

$$F = \frac{4}{3}(\delta x) [2(3\delta x)^4 - (2\delta x)^4 + 2(\delta x)^4] = \frac{592}{3}(\delta x)^5.$$

Whence, from equation 4.9.4:

$$R_N \leq \frac{1}{4!} f_{\max}^{(4)} \left(\frac{1024}{5} - \frac{592}{3}\right) (\delta x)^5 = \frac{28}{90} (\delta x)^5 f_{\max}^{(4)} \quad \dots (4.9.6)$$

so that the potential error is 28 times that of the Simpson's rule formula.

REFERENCES

- (1) BICKLEY, W. G., *J. Math. Phys.*, 27 (1948) 183
- (2) JAHNKE-EMDE, 'Tafeln höherer Funktionen.' (4th edn.) p. 268. Teubner, Leipzig (1948)
- (3) MILNE-THOMSON, L. M. 'The Calculus of finite differences,' (1st edn.), p. 177. Macmillan, London (1951)
- (4) HOBSON, E. W. 'The Theory of Spherical and Ellipsoidal Harmonics,' p. 18, Cambridge (1931)
- (5) — *ibid.*, p. 36
- (6) LOWAN, A. N., DAVIDS, N., and LEVENSON, A. *Bull. Amer. math. Soc.*, 48 (1942), 739
- (7) COURANT, R. and HILBERT, D., 'Methods of Mathematical Physics,' vol. 1 p. 94, Interscience (1953)
- (8) — *ibid.*, vol. 1, p. 91
- (9) BURNETT, D., *Proc. Camb. Phil. Soc.*, xxxiii (1937) 359
- REIZ, A., *Ark. Mat. Astr. Fys.*, 29 iv (1943)
- SALZER, H. E., and ZUCKER, R., *Bull. Amer. math. Soc.*, 55 (1949) 1004
- (10) LANCZOS, C., 'Tables of Chebyshev Polynomials,' p. xiii. *U.S. nat. Bur. Stand. A.M.S. No. 9*
- (11) KOPAL, Z., *Astrophys. J.*, 104 (1946) 74
- (12) FILON, L. N. G., *Proc. Roy. Soc. Edin.*, XLIX (1928) 38
- (13) MILNE, W. E., 'Numerical Calculus,' p. 108. Princeton University Press (1949)

THE SUMMATION OF SERIES

5.1 SOME GENERAL OBSERVATIONS

ALTHOUGH the subject of this chapter might be thought properly to belong to algebra, it is unfortunately true that many recent books on this subject give little practical guidance on appropriate methods of summing series numerically. It is therefore proposed, in this short chapter, to give one or two applications of finite difference operator calculus to this problem.

In the first place we may observe that many series have terms which proceed in unit steps of some argument, for example:

$$\phi(x) = \sum_{r=0}^n \psi(r, x) \quad \dots (5.1.1)$$

where r takes the values $0, 1, 2 \dots n$.

An alternative method of writing equation 5.1.1 makes use of the operator E and is:

$$\phi(x) = \left(\sum_{r=0}^n E^r \right) \psi(0, x) \quad \dots (5.1.2)$$

In this form the 'sum' may readily be found from the well-known formula for the geometric progression, thus:

$$\phi(x) = \frac{E^{n+1} - 1}{E - 1} \psi(0, x) \quad \dots (5.1.3)$$

Suppose now that $\psi(r, x)$ can be represented as the first difference of some new function $\xi(r, x)$ so that

$$\psi(r, x) = \Delta \xi(r, x)$$

equation 5.1.3 now becomes:

$$\phi(x) = \frac{E^{n+1} - 1}{E - 1} \Delta \xi(0, x)$$

or, observing that $E - 1 \equiv \Delta$,

$$\phi(x) = (E^{n+1} - 1) \xi(0, x)$$

whence:

$$\phi(x) = \xi(n+1, x) - \xi(0, x) \quad \dots (5.1.4)$$

which has thus summed the original series—at least formally. It must, unfortunately, be remarked here that only a limited number of functions exist which are recognizable as first differences of other recognizable functions. We shall consider two classes of these in the next section.

5.2 DIFFERENCE FUNCTIONS

An elementary example of a function which can be expressed as the difference of values of a known function has already appeared in section 3.2, equations 3.2.7 and 3.2.8. Thus the factorial function $x^{(m)}$ can always be written as a first difference:

$$x^{(m)} = \Delta \frac{x^{(m+1)}}{(m+1)} \quad \dots (5.2.1)$$

and, since any polynomial in x can be expressed as a polynomial in factorial functions, it follows that any series whose terms are polynomials in x can also be summed.

For example:

$$\begin{aligned} \sum_{z=0}^n z^3 &= \sum_{z=0}^n z(z-1)(z-2) + 3z(z-1) + z \\ &= \sum_{z=0}^n z^{(3)} + 3z^{(2)} + z^{(1)} \\ &= \frac{(n+1)^{(4)}}{4} + (n+1)^{(3)} + \frac{(n+1)^{(2)}}{2} \end{aligned}$$

by virtue of equations 5.1.4 and 5.2.1 and the fact that $0^{(r)} = 0$. Or, by expanding and simplifying:

$$\sum_{z=0}^n z^3 = \left[\frac{n(n+1)}{2} \right]^2$$

which is a well-known result.

Again:

$$\Delta(a^x) = (a^{(x+\delta x)} - a^x) = (a^{\delta x} - 1)a^x \quad \dots (5.2.2)$$

so that:

$$a^x = \Delta \left(\frac{a^x}{a^{\delta x} - 1} \right) \quad \dots (5.2.3)$$

a relation which enables any series which can be put into the form $\sum ca^x$ to be summed.

5.3 THE EULER TRANSFORMATION

Although the actual result of equation 5.1.4 is of limited utility, the analysis which leads to it can be applied to the summation of slowly convergent series with a considerable improvement in rate of convergence.

Thus, consider the series of alternating terms:

$$S_n = u_0 - u_1 + u_2 - u_3 + \dots + (-1)^{n-1} u_{n-1} \dots (5.3.1)$$

This may, by the preceding analysis, be written:

$$S_n = \frac{1+E^n}{1+E} u_0 \dots (5.3.2)$$

Now if $u_n \rightarrow 0$ as $n \rightarrow \infty$ we may assume that:

$$S = \lim_{n \rightarrow \infty} S_n = \frac{u_0}{1+E}$$

and since $E = 1 + \Delta$

$$S = \frac{1}{2} \frac{u_0}{1 + \frac{\Delta}{2}}$$

whence, expanding:

$$S = \frac{1}{2} [u_0 - \frac{1}{2}\Delta u_0 + \frac{1}{4}\Delta^2 u_0 - \frac{1}{8}\Delta^3 u_0 + \dots] \dots (5.3.3)$$

which is usually known as the Euler transformation.

As an instance of the power of this transformation we may mention the example given by BROMWICH⁽¹⁾ who considers the sum:

$$S_1 = 1 - \frac{1}{\sqrt{2}} + \frac{1}{\sqrt{3}} - \frac{1}{\sqrt{4}} + \frac{1}{\sqrt{5}} - \dots$$

and shows that, by taking the first 6 terms of this series and evaluating the differences from the next 7 terms (*i.e.* up to and including $1/\sqrt{12}$), the sum is obtained as $S_1 = 0.6049$ correct to 4 decimal places. To obtain the same accuracy over 10^8 terms of the original series would be needed! Practical points in the use of the Euler transformation are that the first few terms of the original series should be added separately and the transformation applied to the remainder. By taking two different starting points, in this manner, a check on the results of the calculation and an estimate of the accuracy can be had.

One of a number of extensions of the Euler transformation, suggested by Tomlinson Fort, is particularly applicable to an oscillatory series in which $u_{n+1}/u_n \approx \beta^{-1}$. ($\beta > 1$)

We put $u_n = (\beta^{-1})^n v_n$ in the original series and thus obtain:

$$S = \lim_{n \rightarrow \infty} S_n = \frac{v_0}{1 + \beta^{-1}E} = \frac{\beta}{1 + \beta} \left(\frac{v_0}{1 + \frac{\Delta}{1 + \beta}} \right)$$

or

$$S = \frac{\beta}{1 + \beta} \left[v_0 - \frac{1}{1 + \beta} \Delta v_0 + \frac{1}{(1 + \beta)^2} \Delta^2 v_0 - \frac{1}{(1 + \beta)^3} \Delta^3 v_0 + \dots \right] \dots (5.3.4)$$

Since the series in v has terms which are all nearly equal, the difference series (5.3.4) converges rapidly.

5.4 APPLICATION OF THE EULER-MACLAURIN FORMULA

It is evident that the Euler-Maclaurin formula, equation 4.4.5, can be rearranged to give:

$$\begin{aligned} \sum_{r=0}^n f(a + r\delta x) &= \frac{1}{(\delta x)} \int_a^{a+n\delta x} f(x) dx + \frac{1}{2} [f(a) + f(a + n\delta x)] \\ &+ \frac{(\delta x)}{12} [f'(a + n\delta x) - f'(a)] - \frac{(\delta x)^3}{720} [f'''(a + n\delta x) - f'''(a)] \\ &+ \frac{(\delta x)^5}{30240} [f^{(5)}(a + n\delta x) - f^{(5)}(a)] - \dots \dots (5.4.1) \end{aligned}$$

or if, as is usual $\delta x = 1$, and $a = 0$

$$\begin{aligned} \sum_{r=0}^n f(r) &= \int_0^n f(x) dx + \frac{1}{2} [f(0) + f(n)] + \frac{1}{12} [f'(n) - f'(0)] \\ &- \frac{1}{720} [f'''(n) - f'''(0)] + \frac{1}{30240} [f^{(5)}(n) - f^{(5)}(0)] - \dots \dots (5.4.2) \end{aligned}$$

If, in addition, $f^{(s)}(n) \rightarrow 0$ as $n \rightarrow \infty$ for all s , the series converges and:

$$\begin{aligned} \sum_{r=0}^{\infty} f(r) &= \int_0^{\infty} f(x) dx + \frac{1}{2} f(0) - \frac{1}{12} f'(0) + \frac{1}{720} f'''(0) - \frac{1}{30240} f^{(5)}(0) \\ &+ \dots \dots (5.4.3) \end{aligned}$$

These formulae are suitable for application to series of slowly convergent positive terms, but are also useful analytically for the algebraic summation of series. Thus, if $f^{(s)}(x)$ is zero for all s greater

than some lower limit, equation 5.4.3 will give the sum of the infinite series $\sum_{r=0}^{\infty} f(r)$ in finite terms.

Any polynomial series can be summed by means of equation 5.4.2. To take the previously discussed example

$$\begin{aligned}\sum_{r=0}^n r^3 &= \int_0^n r^3 dr + \frac{1}{2}n^3 + \frac{1}{12} \cdot 3n^2 - \frac{1}{720} (6 - 6) \\ &= \frac{n^4}{4} + \frac{n^3}{2} + \frac{n^2}{4} = \left[\frac{n(n+1)}{2} \right]^2\end{aligned}$$

in agreement with the result given in section 5.2.

REFERENCE

- ⁽¹⁾ BROMWICH, T. J. I'A, 'An Introduction to the Theory of Infinite Series,' p. 57, Macmillan, London (1908)

THE SOLUTION OF ORDINARY DIFFERENTIAL EQUATIONS

6.1 INTRODUCTION

THE numerical solution of a differential equation, ordinary or partial, is a problem which has engaged the attention of mathematicians for many years. Numerous methods have been proposed, and some have actually received practical trial; on the whole, however, there appears to be no authoritative statement as to a *best* method, or even a set of such recommended methods.

Upon examining the literature we find, for example, that WHITTAKER and ROBINSON⁽¹⁾ give only one method of solution (that of Bashforth and Adams) and state that this is the best. On the other hand, HARTREE⁽²⁾ makes no mention of this method at all, and a like remark is true of MILNE⁽³⁾, although, in a more recent publication⁽⁴⁾, he points out that most methods are in the nature of modifications of the Bashforth-Adams process.

Ordinary differential equations may be divided roughly into two types for the purpose of numerical solution, depending upon the form of the boundary conditions to be satisfied by solution. The first class of boundary condition may be termed 'one-point,' by which it is to be understood that *all* of the conditions have to be satisfied at a particular point (x, y), say. The second class of boundary condition is a 'two-point' one in which, for example, the function value may be given at one end of the range of integration, and the derivative may be given at another. Of course, with equations of order higher than the second 'multi-point' boundary conditions are possible; we shall not, however, be concerned with these in the present work. It may be mentioned that the one-point condition is sometimes said to lead to a 'marching' problem, and the two-point condition to a 'jury' problem, for reasons which will become clear later.

In this chapter we shall be concerned chiefly with one-point problems and shall mention the more complicated two-point variant only in so far as it can be solved by essentially one-point methods. A more satisfactory approach will be given later, in Chapter 8, section 8.6, when relaxation solutions are discussed.

6.2 INITIAL VALUES

In many of the finite difference methods of solution it is necessary to have one or more values of the solution, and possibly of its derivatives, near to the starting point. These are usually best obtained by a process rather different from that used in the remainder of the solution and we shall give, briefly, the two methods most frequently used.

The first method is that of Picard, and depends upon the 'guessing' of an initial solution. Consider the first order equation:

$$\frac{dy}{dx} = f(x, y) \quad \dots (6.2.1)$$

and assume that near to the boundary, (a, b) say, we can guess an approximation:

$$y = g_1(x) \quad \dots (6.2.2)$$

We then substitute in equation 6.2.1 and integrate to obtain:

$$y = b + \int_a^x f[x, g_1(x)] \quad [= g_2(x), \text{ say}]$$

the process is now repeated:

$$y = b + \int_a^x f[x, g_2(x)] = g_3(x) \quad \text{and so on.}$$

This process, if it converges, will give an approximation of any desired degree of accuracy near to the initial point.

The range of application of this method is fairly limited, since the approximating functions $f[x, g_r(x)]$ have to be integrable in closed form. As an example of the use of the method consider the equation: $\frac{dy}{dx} = x^2 + 2xy$ with $x = 0$ when $y = 0$. Take:

$$g_1(x) = \frac{x^3}{3}$$

then

$$g_2(x) = \int_0^x \left(x^2 + \frac{2x^4}{3} \right) dx = \frac{x^3}{3} + \frac{2x^5}{3.5}$$

$$g_3(x) = \int_0^x \left(x^2 + \frac{2x^4}{3} + \frac{4x^6}{3.5} \right) dx = \frac{x^3}{3} + \frac{2x^5}{3.5} + \frac{4x^7}{3.5.7}$$

$$g_4(x) = \int_0^x \left(x^2 + \frac{2x^4}{3} + \frac{4x^6}{3.5} + \frac{8x^8}{3.5.7} \right) dx \\ = \frac{x^3}{3} + \frac{2x^5}{3.5} + \frac{4x^7}{3.5.7} + \frac{8x^9}{3.5.7.9} \quad \text{etc.}$$

INITIAL VALUES

The reader will easily convince himself that the expression agrees with the analytic solution:

$$y = -\frac{1}{2}x + \frac{1}{2}e^{x^2} \int_0^x e^{-x^2} dx$$

when the latter is expanded in a series. It is also worth noticing that the Picard solution, in the example, is much easier to compute than that derived from the analytic solution.

The second method of obtaining an approximate initial solution is by means of a Taylor series expansion. Thus for $y = b$ at $x = a$ we have:

$$y = b + \frac{(x-a)}{1!} g'(a) + \frac{(x-a)^2}{2!} g''(a) + \dots + \frac{(x-a)^n}{n!} g^n(a) + \dots \quad \dots (6.2.3)$$

and, from the differential equation:

$$\frac{dy}{dx} = g'(x) = f(x, y) \quad \dots (6.2.4)$$

$$\text{giving:} \quad g'(a) = f(a, b) \quad \dots (6.2.5)$$

Differentiating 6.2.4 and substituting the boundary condition and 6.2.5:

$$\frac{d^2y}{dx^2} = \frac{\partial f(x, y)}{\partial x} + \frac{\partial f(x, y)}{\partial y} \cdot \frac{dy}{dx}$$

$$\text{i.e.} \quad g''(a) = \left[\frac{\partial f(x, y)}{\partial x} \right]_{x=a, y=b} + \left[\frac{\partial f(x, y)}{\partial y} \right]_{x=a, y=b} f(a, b).$$

and similarly for derivatives of higher order.

The process may be illustrated by our previous example. Since $x = 0, y = 0$, the appropriate expansion (really a Maclaurin series) is:

$$y = xg'(0) + \frac{x^2}{2!} g''(0) + \frac{x^3}{3!} g^{(3)}(0) + \dots$$

The differential equation, together with $x = 0, y = 0$ gives:

$$\begin{aligned} g'(0) &= [x^2 + 2xy]_{0,0} = 0 \\ g''(0) &= [2x + 2y + 2xy']_{0,0} = 0 \\ g^{(3)}(0) &= [2 + 2y' + 2y'' + 2xy'']_{0,0} = 2 \\ g^{(4)}(0) &= [4y'' + 2y''' + 2xy''']_{0,0} = 0 \\ g^{(5)}(0) &= [6y^{(3)} + 2y^{(4)} + 2xy^{(4)}]_{0,0} = 16 \\ g^{(6)}(0) &= [8y^{(4)} + 2y^{(5)} + 2xy^{(5)}]_{0,0} = 0 \\ g^{(7)}(0) &= [10y^{(5)} + 2y^{(6)} + 2xy^{(6)}]_{0,0} = 12.16 \quad \text{etc.} \end{aligned}$$

which will be seen to agree with the series previously deduced.

It will be noticed that the Taylor series method, in this example, involves slightly more work than the Picard technique; on the other hand, differentiations involved can *always* be performed, for the commonly encountered functions at least, whereas the integrations of the latter method can only be carried out for relatively simple functions.

On the other hand, it is possible to carry out the Picard integrations numerically, in the required range, and this is often a simple and rapid procedure. Certain integration formulae are required for this operation, and these are obtainable by the same technique as was used in section 4.3. A useful set of results is the following:

$$\begin{aligned}\int_0^{(\delta x)} f(x) dx &= \frac{(\delta x)}{720} [251f(0) + 646f(1) - 264f(2) + 106f(3) - 19f(4)] \\ &\quad + \frac{27(\delta x)^6}{1440} f_m^{(6)} \\ \int_0^{2(\delta x)} f(x) dx &= \frac{(\delta x)}{90} [29f(0) + 124f(1) + 24f(2) + 4f(3) - f(4)] \\ &\quad + \frac{16(\delta x)^6}{1440} f_m^{(6)} \\ \int_0^{3(\delta x)} f(x) dx &= \frac{(\delta x)}{80} [27f(0) + 102f(1) + 72f(2) + 42f(3) - 3f(4)] \\ &\quad + \frac{27(\delta x)^6}{1440} f_m^{(6)} \\ \int_0^{4(\delta x)} f(x) dx &= \frac{4(\delta x)}{90} [7f(0) + 32f(1) + 12f(2) + 32f(3) + 7f(4)] \\ &\quad - \frac{8(\delta x)^7}{945} f_m^{(7)} \\ &\quad \dots 6.2.6\end{aligned}$$

where $f_m^{(n)}$ is some value of $f^{(n)}$ in the interval $(0, 4)$. These are readily obtained by integrating the approximating polynomial:

$$\begin{aligned}f_f(x) &= f(0) + x\Delta f(0) + \frac{x(x-1)}{2!} \Delta^2 f(0) + \frac{x(x-1)(x-2)}{3!} \Delta^3 f(0) \\ &\quad + \frac{x(x-1)(x-2)(x-3)}{4!} \Delta^4 f(0) \quad \dots (6.2.7)\end{aligned}$$

between the appropriate limits, and evaluating the error terms in the manner described in section 4.8. They form one of a set of similar results obtainable from approximating polynomials of

various degrees, and are given here, in preference to other members of the group, because they are a compromise between accuracy and ease of working. In suitable cases equations 6.2.6, in conjunction with the Picard technique, will give initial values at $(\delta x) = 0.2$ to an accuracy of slightly better than one in 10^5 .

An even more satisfactory technique is to make use of values of $f(x)$ on either side of the origin. This is particularly useful when central difference methods are to be used at a later stage to carry on the solution. The relevant formulae can be obtained, either by integrating equation 6.2.7, or from Everett's formula. For the five points $(-2, -1, 0, +1, +2)$ the results are:

$$\begin{aligned}\int_{-2(\delta x)}^0 f(x) dx &= \frac{(\delta x)}{90} [29f(-2) + 124f(-1) + 24f(0) + 4f(1) - f(2)] \\ &\quad - \frac{41(\delta x)^6}{3600} f_m^{(6)} \\ \int_{-(\delta x)}^0 f(x) dx &= \frac{(\delta x)}{720} [-19f(-2) + 346f(-1) + 456f(0) - 74f(1) + 11f(2)] \\ &\quad - \frac{11(\delta x)^6}{1440} f_m^{(6)} \\ \int_0^{(\delta x)} f(x) dx &= \frac{(\delta x)}{720} [11f(-2) - 74f(-1) + 456f(0) + 346f(1) - 19f(2)] \\ &\quad + \frac{11(\delta x)^6}{1440} f_m^{(6)} \\ \int_0^{2(\delta x)} f(x) dx &= \frac{(\delta x)}{90} [-f(-2) + 4f(-1) + 24f(0) + 124f(1) + 29f(2)] \\ &\quad + \frac{41(\delta x)^6}{3600} f_m^{(6)} \\ &\quad \dots (6.2.8)\end{aligned}$$

Although we have given the Picard method and certain integration formulae for its application, it is not our practice to use this technique unless the direct analytic integration is possible. We have found that the Taylor series approach is more easily applicable and is, furthermore, readily extended to equations of higher degree than the first. To take only one example, consider:

$$y'' = f(x, y, y') \quad \dots (6.2.9)$$

with initial conditions $x = x_0, y = y_0, y' = y'_0$.

Equation 6.2.9 gives:

$$y_0'' = f(x_0, y_0, y_0')$$

by direct substitution of the initial conditions. Differentiating:

$$y^{(3)} = \frac{d}{dx} f(x, y, y') = \frac{\partial f}{\partial x} + \frac{\partial f}{\partial y} \frac{dy}{dx} + \frac{\partial f}{\partial y'} \frac{d^2 y}{dx^2}$$

so that:

$$y_0^{(3)} = \left(\frac{\partial f}{\partial x} \right)_{x=x_0} + \left(\frac{\partial f}{\partial y} \right)_0 y_0' + \left(\frac{\partial f}{\partial y'} \right)_0 y_0''$$

and so on.

6.3 EQUATIONS OF THE FIRST ORDER

The 'classical' method of solution of a first order differential equation:

$$\frac{dy}{dx} = f(x, y) \quad \dots (6.3.1)$$

subject to one-point boundary conditions, is that of Runge. It is based upon an idea originally due to Euler and has been extended by Kutta, Heun, and Piaggio. The method is normally used to obtain an initial solution with which to start one of the Bashforth and Adams type methods. Since, in our opinion, it has no advantages over the Taylor series method or the numerical version of the Picard process, we shall do no more than state the relevant formulae in the two most useful cases.

The first Runge-Kutta approximation is the equivalent of Simpson's rule applied to the integration of the $f(x, y)$. If, in equation 6.3.1, we know that $y = y_0$ when $x = x_0$ and require y_1 , the value of y when $x = x_0 + (\delta x)$. Then:

$$y_1 = y_0 + \frac{1}{6}(k_1 + 4k_2 + k_3) \quad \dots (6.3.2)$$

where

$$k_1 = (\delta x) \cdot f(x_0, y_0)$$

$$k_2 = (\delta x) \cdot f(x_0 + \frac{1}{2}(\delta x), y_0 + \frac{1}{2}k_1)$$

$$k_3 = (\delta x) \cdot f(x_0 + (\delta x), y_0 + 2k_2 - k_1)$$

the error in y_1 being of order $(\delta x)^4$.

If an error of order $(\delta x)^5$ is required, the Kutta fourth order process:

$$y_1 = y_0 + \frac{1}{6}(k_1 + 2k_2 + 2k_3 + k_4) \quad \dots (6.3.3)$$

where:

$$k_1 = (\delta x) f(x_0, y_0)$$

$$k_2 = (\delta x) f(x_0 + \frac{1}{2}(\delta x), y_0 + \frac{1}{2}k_1)$$

$$k_3 = (\delta x) f(x_0 + \frac{1}{2}(\delta x), y_0 + \frac{1}{2}k_2)$$

$$k_4 = (\delta x) f(x_0 + (\delta x), y_0 + k_3) \quad \text{may be used.}$$

By increasing the number of steps, still higher accuracy may be attained, but the formulae are considerably more complicated and have not been much used. For a more extensive treatment of Runge-Kutta formulae the reader is referred to the standard literature.^(6, 6)

A major criticism of the Runge-Kutta method of solution lies in the fact that it contains, in itself, no means of checking either the accuracy of the solution obtained, or the actual arithmetic by means of which such a solution is derived. The objection is removed in the Bashforth-Adams process for continuing the solution, which we shall now describe.

We first notice that the method initially uses the values of the function previously computed to obtain new values and then uses the new values, and the equations, to check and correct the result. The differential equation itself is used as a source of values of the derivative. We shall require an extrapolation formula which uses only differences which can be computed from existing function values, and this implies the use of backward differences (see section 3.1). Now:

$$f(x) = E^x f(0)$$

$$\text{and} \quad \nabla \equiv (E - 1)/E$$

$$\text{whence} \quad f(x) = f(0)/(1 - \nabla)^x$$

$$= \left[1 + x\nabla + \frac{x(x+1)}{2!} \nabla^2 + \frac{x(x+1)(x+2)}{3!} \nabla^3 + \dots \right] f(0) \quad \dots (6.3.4)$$

Now

$$\int_0^1 f(x) dx = \left[1 + \frac{1}{2} \nabla + \frac{5}{12} \nabla^2 + \frac{3}{8} \nabla^3 + \frac{251}{720} \nabla^4 + \frac{95}{288} \nabla^5 + \frac{19087}{60480} \nabla^6 + \dots \right] f(0)$$

or, for interval (δx) :

$$\int_0^{(\delta x)} f(x) dx = (\delta x) \left(1 + \frac{1}{2} \nabla + \frac{5}{12} \nabla^2 + \frac{3}{8} \nabla^3 + \frac{251}{720} \nabla^4 + \dots \right) f(0) \quad \dots (6.3.5)$$

To solve the equation:

$$\frac{dy}{dx} = f(x, y) \text{ or } y_1 = y_0 + \int_0^{\delta x} f(x, y) dx$$

we assume that five initial values, (x_{-4}, y_{-4}) , (x_{-3}, y_{-3}) , (x_{-2}, y_{-2}) , (x_{-1}, y_{-1}) , (x_0, y_0) , say, have been calculated by one of the methods previously described. From these values the table

x	y	$dy/dx = f(x, y)$	∇f	$\nabla^2 f$	$\nabla^3 f$	$\nabla^4 f$	$\nabla^5 f$
x_{-4}	y_{-4}	f_{-4}					
x_{-3}	y_{-3}	f_{-3}	∇f_{-3}				
x_{-2}	y_{-2}	f_{-2}	∇f_{-2}	$\nabla^2 f_{-2}$			
x_{-1}	y_{-1}	f_{-1}	∇f_{-1}	$\nabla^2 f_{-1}$	$\nabla^3 f_{-1}$		
x_0	y_0	f_0	∇f_0	$\nabla^2 f_0$	$\nabla^3 f_0$	$\nabla^4 f_0$	
<hr style="border-top: 1px dashed black;"/>							
x_1	y_1	f_1	∇f_1	$\nabla^2 f_1$	$\nabla^3 f_1$	$\nabla^4 f_1$	$\nabla^5 f_1$

is constructed, as far as the broken line ----. Using the values of $\nabla^r f_0$ up to $r = 4$ in the integration formula (6.3.5) the value y_1 may be obtained; this enables f_1 to be calculated and, once this is known, the differences $\nabla^r f_1$ can be evaluated.

To check the accuracy of the prediction we may form the integral:

$$\int_0^{-(\delta x)} f(x) = -(\delta x) \left(1 - \frac{1}{2}\nabla - \frac{1}{12}\nabla^2 - \frac{1}{24}\nabla^3 - \frac{1}{720}\nabla^4 - \frac{3}{160}\nabla^5 - \frac{883}{60480}\nabla^6 - \dots \right) f(1)$$

or, in terms of y_1, y_0 and the differences below the broken line which can now be calculated:

$$y_1 = y_0 + (\delta x) \left(f_1 - \frac{1}{2}\nabla f_1 - \frac{1}{12}\nabla^2 f_1 - \frac{1}{24}\nabla^3 f_1 - \frac{1}{720}\nabla^4 f_1 - \frac{3}{160}\nabla^5 f_1 - \frac{883}{60480}\nabla^6 f_1 - \dots \right) \quad (6.3.6)$$

This formula, up to and including ∇^4 , is over ten times more accurate than equation 6.3.5, so that, if the two values of y_1 do not differ by more than 5 units in the last decimal place required, we may accept the value derived from equation 6.3.6 as correct. If the deviation is greater than this, the *new* value of y_1 can be used to recompute $f_1, \nabla f_1$, etc. and equation 6.3.6 re-applied until agreement between consecutive values is reached.

The process is now repeated, using subscript 1 instead of 0 as previously, and y_2 is thus derived.

A modification of the Bashforth-Adams method, due to Milne will now be described. The central idea is again to predict the value of y_{n+1} from those of y_n and earlier ordinates, and then to correct it (if necessary) by means of a comparison formula which will, in general, involve y_{n+1} itself. Milne recommends that the open Newton-Cotes formula, given by equation 4.3.10, be used to 'predict' and Simpson's rule (equation 4.3.4) to correct. These formulae become:

'Predictor'

$$y_{n+1} = y_{n-3} + \frac{4}{3}(\delta x)(2f_n - f_{n-1} + 2f_{n-2}) + \frac{28}{90}(\delta x)^5 f_m^{(6)} \quad \dots (6.3.7)$$

'Corrector'

$$y_{n+1} = y_{n-1} + \frac{(\delta x)}{3}(f_{n+1} + 4f_n + f_{n-1}) - \frac{1}{90}(\delta x)^5 f_m^{(6)} \quad \dots (6.3.8)$$

in terms of the solution values y_n and the values f in equation 6.3.1. The interval, (δx) , should be so chosen that the two values of y_{n+1} differ by less than 14 in the last decimal place required. If this is done, a comparison of the error terms in equations 6.3.7 and 6.3.8 shows that the value derived from 6.3.8 can be taken as correct within the number of places required.

We conclude this section by describing a procedure which has been described by MILNE⁽⁷⁾ and, in modified form, by HARTREE⁽⁸⁾.

The equation is, as usual, 6.3.1, and from it we can obtain:

$$y' = f$$

$$y'' = \frac{d^2 y}{dx^2} = \frac{\partial f}{\partial x} + f \frac{\partial f}{\partial y} \quad \dots (6.3.9)$$

Now the Taylor series expansion, based on y_n , gives:

$$y_{n+1} = y_n + (\delta x)y'_n + \frac{(\delta x)^2}{2!}y''_n + \frac{(\delta x)^3}{3!}y'''_n + \frac{(\delta x)^4}{4!}y^{(4)}_n + \dots (6.3.10)$$

whence:

$$-\frac{1}{2}(\delta x)y'_{n+1} = -\frac{1}{2}(\delta x)y'_n - \frac{1}{2}(\delta x)^2y''_n - \frac{1}{2} \cdot \frac{(\delta x)^3}{2!}y'''_n - \frac{1}{2} \cdot \frac{(\delta x)^4}{3!}y^{(4)}_n + \dots$$

$$+ \frac{1}{12}(\delta x)^2(y'')_{n+1} = \frac{1}{12}(\delta x)^2y''_n + \frac{1}{12}(\delta x)^3y'''_n + \frac{1}{12} \cdot \frac{(\delta x)^4}{2!}y^{(4)}_n + \dots$$

or, adding

$$y_{n+1} = y_n + \frac{1}{2}(\delta x)(y'_{n+1} + y'_n) - \frac{1}{12}(\delta x)^2(y''_{n+1} - y''_n) + R_e \quad \dots(6.3.11)$$

where R_e can be shown to be given by:

$$R_e = (\delta x)^5 y^{(5)}_n / 720.$$

Again, by forming the sum of:

$$\begin{aligned} y_{n+1} &= y_n + (\delta x)y'_n + \frac{(\delta x)^2}{2!}y''_n + \frac{(\delta x)^3}{3!}y^{(3)}_n + \frac{(\delta x)^4}{4!}y^{(4)}_n + \dots \\ -2y_n &= -2y_{n-1} - 2(\delta x)y'_{n-1} - 2\frac{(\delta x)^2}{2!}y''_{n-1} - 2\frac{(\delta x)^3}{3!}y^{(3)}_{n-1} \\ &\quad - 2\frac{(\delta x)^4}{4!}y^{(4)}_{n-1} - \dots \\ -y_{n-2} &= -y_{n-1} + (\delta x)y'_{n-1} - \frac{(\delta x)^2}{2!}y''_{n-1} + \frac{(\delta x)^3}{3!}y^{(3)}_{n-1} \\ &\quad - \frac{(\delta x)^4}{4!}y^{(4)}_{n-1} + \dots \\ -\frac{1}{2}(\delta x)^2 y''_n &= -\frac{(\delta x)^2}{2}y''_{n-1} - \frac{(\delta x)^3}{2}y^{(3)}_{n-1} \\ &\quad - \frac{(\delta x)^4}{4}y^{(4)}_{n-1} - \dots \end{aligned}$$

and then substituting for the derivatives $y'_n, y''_n, y^{(3)}_n$ in terms of the Taylor series based on y_{n-1} , it is easily seen that:

$$y_{n+1} = y_{n-2} + 3(y_n - y_{n-1}) + (\delta x)^2(y''_n - y''_{n-1}) + R_p \quad \dots(6.3.12)$$

where

$$R_p = 60(\delta x)^5 y^{(5)}_n / 720.$$

Formula 6.3.12 can be used to predict the value of y_{n+1} given those of y_n, y_{n-1} and y_{n-2} , the differential equation itself being used to determine y''_n and y''_{n-1} (from equation 6.3.9).

From these predicted values a corrected value can be calculated via equation 6.3.11. If the difference between the predicted and corrected values is less than 30 in the last decimal place required, the latter is taken to be correct. If the difference exceeds this figure the corrected value is used to recalculate the derivatives y'_{n+1} ,

y''_{n+1} in equation 6.3.11, and a new and more accurate value of y_{n+1} is thus obtained. This process is continued until two successive values of y_{n+1} do not differ within the limits of accuracy required. It should be noted that, with a properly chosen value of (δx) , a re-application of equation 6.3.11 should not often be required. If it is, there is a clear indication that a change in (δx) is needed. In a similar manner, if repeated use of 6.3.11 is never required, it suggests that too small a value of (δx) is being used. A working rule is that one re-application of equation 6.3.11 should be required every five or so steps.

6.4 EQUATIONS OF THE SECOND ORDER

The canonical form of the differential equation of the second order is:

$$\frac{d^2y}{dx^2} = f(x, y) \quad \dots(6.4.1)$$

and any equation of the second order, linear in the derivatives, can be reduced to the form of equation 6.4.1. Thus:

$$\frac{d^2y}{dx^2} + a(x)\frac{dy}{dx} + b(x, y) = 0$$

reduces to 6.4.1 when the transformation:

$$y = Y \exp \left[-\frac{1}{2} \int a(x) dx \right] \quad \dots(6.4.2)$$

is applied. The transformed equation is:

$$Y'' = \left[\frac{1}{2}a'(x) + \frac{1}{4}\{a(x)\}^2 \right] Y - b(x, Y) \exp \frac{1}{2} \int a(x) dx \quad \dots(6.4.3)$$

but unless $b(x, y)$ is of the form $y \cdot c(x)$, this is not likely to be suitable for numerical calculation because of the exponential factor. Furthermore, the transformation destroys any periodicity which may exist in solutions of the original equation.

By means of the transformation:

$$z = y' \quad \dots(6.4.4)$$

the general second order equation:

$$y'' = f(x, y, y') \quad \dots(6.4.5)$$

is reduced to the pair of simultaneous first order equations:

$$\left. \begin{aligned} z' &= f(x, y, z) \\ z &= y' \end{aligned} \right\} \quad \dots(6.4.6)$$

which may be solved by the methods discussed in section 6.5. This reduction has the merit of simplicity, but is not to be recommended for equations in which y' is absent.

The Bashforth-Adams method can be applied directly to equation 6.4.5 when a suitable set of initial values has been obtained; the detailed procedure is as follows:

- (1) Assume that y_n'' , y_{n-1}'' . . . *etc.* are known.
- (2) Find y_{n+1}' from y_n'' , y_{n-1}'' *etc.* by means of 6.3.5.
- (3) Assume that y_n' , y_{n-1}' *etc.* are known.
- (4) Find y_{n+1} from y_{n+1}' , y_n' . . . *etc.* by means of 6.3.6.
- (5) Using these values of y_{n+1}' , y_{n+1} compute y_{n+1}'' from 6.4.5.
- (6) Recompute y_{n+1}' by means of 6.3.6.

At this stage the usual error estimation can be made and, if necessary, one or more extra cycles of the process can be used to obtain the desired accuracy.

The Bashforth-Adams backward difference formulae are not the only ones which can be used in this process. Thus MILNE⁽⁹⁾ suggests that the open Newton-Cotes formula (4.3.10), in the form:

$$y_{n+1}' = y_{n-3}' + \frac{2}{3}(\delta x)(2y_n'' - y_{n-1}'' + 2y_{n-2}'')$$

should be used in step (2) (*supra*), and that this should be followed by Simpson's rule:

$$y_{n+1} = y_{n-1} + \frac{(\delta x)}{3} (y_{n+1}' + 4y_n' + y_{n-1}')$$

in step (4) and again in step (6), this time as:

$$y_{n+1}' = y_{n-1}' + \frac{(\delta x)}{3} (y_{n+1}'' + 4y_n'' + y_{n-1}'').$$

For second order equations in which the first derivative is absent (*i.e.* the canonical form 6.4.1), a simplified procedure is available which, in effect, uses a double integration formula to pass directly from y'' to y .

One such formula has been obtained in equation 4.5.3; it may be rewritten, for the present purpose:

$$y_{n+1} = 2y_n - y_{n-1} + (\delta x)^2 (y_n'' + \frac{1}{12}\delta^2 y_n'' - \frac{1}{240}\delta^4 y_n'' + \dots) \quad \dots (6.4.7)$$

or

$$y_{n+1} = 2y_n - y_{n-1} + \frac{(\delta x)^2}{12} (y_{n+1}'' + 10y_n'' + y_{n-1}'') - \frac{1}{240}(\delta x)^6 y_m^{(6)} \quad \dots (6.4.7a)$$

This result may be used as a corrector and the less accurate:

$$y_{n+1} = y_n + y_{n-2} - y_{n-3} + \frac{(\delta x)^2}{4} (5y_n'' + 2y_{n-1}'' + 5y_{n-2}'') + \frac{17}{240}(\delta x)^6 y_m^{(6)} \quad \dots (6.4.8)$$

as a predictor. As an alternative to equation 6.4.8:

$$y_{n+1} = 2y_{n-1} - y_{n-3} + 4(\delta x)^2 (y_{n-1}'' + \frac{1}{3}\delta^2 y_{n-1}'') + \frac{19}{240}(\delta x)^6 y_m^{(6)} \quad \dots (6.4.9)$$

which may be regarded as the analogue of equation 6.4.7, but for double the interval, is perhaps simpler to apply.

Another, and more recent, method for solving the canonical form of the second-order equation is that of DE VOGELAERE⁽¹⁰⁾.

The equations used are:

$$y_{n+\frac{1}{2}} = y_n + \frac{1}{2}(\delta x)y_n' + \frac{1}{24}(\delta x)^2(4y_n'' - y_{n-\frac{1}{2}}'') + 0(\delta x^4)$$

$$y_{n+1} = y_n + (\delta x)y_n' + \frac{1}{6}(\delta x)^2(y_n'' + 2y_{n+\frac{1}{2}}'') + 0(\delta x^5)$$

$$y_{n+\frac{1}{2}}' = y_n' + \frac{1}{6}(\delta x)(y_n'' + 4y_{n+\frac{1}{2}}'' + y_{n+1}'') + 0(\delta x^5)$$

which, as in the Runge-Kutta process, allow the progressive solution of the equation. The advantage of de Vogelaere's method is that it requires the evaluation of $y_n'' (= f(x_n, y_n))$ only twice per step.

To start the solution it is necessary to know y_0 and y_0' and $y_{-\frac{1}{2}}'$. It is usual to compute the latter from:

$$y_{-\frac{1}{2}} = y_0 - \frac{1}{2}(\delta x)y_0' + \frac{1}{8}(\delta x)^2 y_0''.$$

When the form of $f(x, y)$ in equation 6.4.1 is such that its calculation for any particular values of x and y is difficult or tedious, a procedure due to Jennings, Fox and Goodwin often enables a large interval by differencing to be used, with consequent reduction in the number of values of $f(x, y)$ to be calculated. In essence, the method uses a large value of (δx) and a simple formula to obtain a rough solution, and then uses the differences obtained from this solution to enable a better approximation to be made.

Thus, suppose that a set of values $y_{(m-1)}$ at interval (δx) have been obtained; then, from equation 4.5.3

$$\delta^2 y = (\delta x)^2 (y'' + \frac{1}{12}\delta^2 y'' - \frac{1}{240}\delta^4 y'' + \dots)$$

and, if we estimate $\delta^2 y''$, $\delta^4 y''$ from the previous approximation,

$$\delta^2 y_{(m)} = (\delta x)^2 [y_{(m)}'' + \frac{1}{12}\delta^2 y_{(m-1)}'' - \frac{1}{240}\delta^4 y_{(m-1)}'' + \dots]$$

$$\text{or } \delta^2 y_{(m)} = (\delta x)^2 [f(x, y_{(m)}) + \frac{1}{12}\delta^2 y_{(m-1)}'' - \frac{1}{240}\delta^4 y_{(m-1)}'' \dots].$$

Thus, if $y_{(m)n+1}$, $y_{(m)n}$, $y_{(m)n-1}$ are consecutive points in the m th approximation:

$$y_{(m)n+1} = 2y_{(m)n} - y_{(m)n-1} + (\delta x)^2 \left[f(x, y_{(m)n}) + \frac{1}{12} \delta^2 y_{(m-1)n}'' - \frac{1}{240} \delta^4 y_{(m-1)n}^{(4)} \dots \right] \quad \dots (6.4.10)$$

a formula which enables the solution to be continued from the n th point of the m th approximation to the $(n+1)$ th point.

The same general method can be applied to equations of the first order and a comparison of such techniques has been made by Fox and GOODWIN⁽¹¹⁾. A particular calculation, in which an economy of labour has been achieved, is the tracing of paraxial rays through an electron lens system where the field distribution is given by complete elliptic integrals which are awkward to manipulate. This problem has been extensively studied by JENNINGS⁽¹²⁾

6.5 SIMULTANEOUS DIFFERENTIAL EQUATIONS

The typical set of simultaneous differential equations of the first order may be written:

$$y_m' = f_m(x, y_1, y_2, \dots, y_n) \quad (m = 1 \dots n). \quad \dots (6.5.1)$$

The methods of solution are almost identical with those suggested for the ordinary first order equations, and are typified by the following scheme.

- (1) Obtain initial values from a Taylor series expansion and the given boundary conditions.
- (2) From $y_{m,0}$, $y_{m,0}$ ($m = 1 \dots n$) predict values of $y_{m,1}$ by means of equation 6.3.5 or 6.3.7.
- (3) Using the values of $y_{m,1}$ ($m = 1 \dots n$) obtained in (2) recalculate $y_{m,1}$ from the differential equations and 6.3.6 or 6.3.8.
- (4) If significant changes occur repeat the process from (3) until adequate accuracy has been obtained.

It should be noticed that this technique can be applied, as indicated in section 6.4, to obtain the solution of a second order equation.

In a similar manner, the set of equations of the second order:

$$y_m'' = f_m(x, y_1, y_2, \dots, y_n) \quad (m = 1 \dots n) \quad \dots (6.5.2)$$

may be solved, using the repeated integration formulae of equations 6.4.8 and 6.4.9 as predictors, and 6.4.7 as a corrector.

6.6 EQUATIONS OF HIGHER ORDER

When such equations, either single or simultaneous, are encountered, the only really practical method of solution appears to be to reduce them to sets of first order simultaneous equations by means of successive transformations of the type in equation 6.4.4. The actual integrations at any step should proceed from the highest order derivative downwards and (6.3.5, 6.3.6) (6.3.7, 6.3.8) or (6.3.11, 6.3.12) are suitable pairs of integration formulae.

6.7 MULTI-POINT BOUNDARY CONDITIONS

Although we shall defer consideration of this type of solution until Chapter 8, section 6, it is appropriate to indicate here the mode of application of the methods just discussed to this problem.

First consider a second order differential equation

$$y'' = f(y', y, x). \quad \dots (6.7.1)$$

If this is to satisfy:

$$y = y_0 \text{ at } x = x_0$$

$$y = y_1 \text{ at } x = x_1$$

there is, of course, insufficient data with which to start a finite difference integration from $x = x_0$. Suppose, however, that 6.7.1 is a linear equation; then if $y = l_1(x)$ and $y = l_2(x)$ are any two solutions, $y = Al_1(x) + Bl_2(x)$ is also a solution.

To solve the given equation all that is needed is the following. From $x = x_0$ start any two integrations having $y = y_0$ but *different* values of y' at $x = x_0$. Carry these solutions up to $x = x_1$ and assume these solutions to be $l_1(x)$ and $l_2(x)$. The required solution is then:

$$y = Al_1(x) + Bl_2(x)$$

where A and B are the solutions of:

$$y_0 = Al_1(x_0) + Bl_2(x_0) \quad (\text{i.e. } A + B = 1)$$

$$y_1 = Al_1(x_1) + Bl_2(x_1).$$

When equation 6.7.1 is not linear this technique will not work, since $Al_1(x) + Bl_2(x)$ is no longer, in general, a solution of the equation. In this event, a possible method of approach is to guess a pair of initial directions which produce solutions straddling the required values at $x = x_1$. Linear interpolation then gives a better initial direction, and the process is repeated until the desired accuracy is attained. Since most of the non-linear equations encountered in practice, have solutions which vary exponentially as some power

of the independent variable, this method is, however, seldom applicable over any considerable range.

An extension of this procedure is to take more than two trial solutions and then to use a non-linear interpolation procedure to obtain the correct initial value.

6.8. EIGENVALUE PROBLEMS

It is often required to obtain the 'eigenvalues' of a differential operator. By this is meant those values of λ_n for which the equation:

$$L\psi_n(x) = \lambda_n\psi_n(x) \quad \dots (6.8.1)$$

has solutions which satisfy specified conditions at $x = a$ and $x = b$. A typical form⁽¹³⁾ for the operator L is:

$$L \equiv g(x) - \frac{d^2}{dx^2} \quad \dots (6.8.2)$$

and (a, b) is often $(0, \infty)$. The trivial solution $\psi_n(x) = 0$ is excluded.

Depending upon the form of L , and of the boundary conditions, λ_n may have discrete values only, or a complete spectrum. The determination of the values of λ_n can be carried out by the process of constructing trial solutions for different assumed values of λ and then following the process outlined in section 6.7 to obtain more accurate values.

It should be emphasized that numerical calculation of eigenvalues should never be used as an alternative to a proper analytical investigation of the equation. In many real physical situations such an investigation will lead to a fruitful result.

An alternative approach to the determination of approximate eigenvalues is *via* the equivalent set of simultaneous difference equations, but this is dealt with more fully in Chapter 8.

REFERENCES

- (1) WHITTAKER, E. T. and ROBINSON, G., 'The Calculus of Observations,' 4th edn. Blackie (1927)
- (2) HARTREE, D. R., 'Numerical Analysis,' Oxford University Press, London (1952)
- (3) MILNE, W. E., 'Numerical Calculus,' Princeton University Press (1949)
- (4) — 'Numerical Solution of Differential Equations,' Wiley, New York (1953)
- (5) LEVY, H. and BAGGOTT, E. A., 'Numerical Studies in Differential Equations,' p. 91 *et seq.* Watts, London (1934)

REFERENCES

- (6) MILNE, W. E., 'Numerical Solution of Differential Equations,' p. 72 *et seq.* Wiley, New York (1953).
- (7) — 'Numerical Solution of Differential Equations,' p. 76 *et seq.* Wiley, New York (1953)
- (8) HARTREE, D. R., 'Numerical Analysis,' p. 133 *et seq.* Oxford University Press, London (1952)
- (9) MILNE, W. E., 'Numerical Calculus,' p. 139, Princeton University Press (1949)
- (10) DE VOGELAERE, R., *J. Res. natn. Bur. Stand.*, 54 (1955) 119
- (11) FOX, L. and GOODWIN, E. T., *Proc. Camb. Phil. Soc.*, 45 (1949) 373
- (12) JENNINGS, J. C. E. and PRATT, R. G., *Proc. Phys. Soc.*, B, LXVIII (1955) 526
- (13) TITCHMARSH, E. C., 'Eigenfunction Expansions,' Oxford University Press, London (1946)

SIMULTANEOUS LINEAR EQUATIONS

7.1 PRELIMINARY REMARKS

ALTHOUGH simultaneous linear equations arise in a straightforward manner in many problems of physical and engineering science, they have attained a much greater importance in recent years because of their application to the solution of various sorts of differential equation—both ordinary and partial.

Whereas the 'classical' application of simultaneous equations, to such problems as an ordnance survey, rarely results in a set of more than 10 equations in 10 unknowns, many of the more modern approaches to differential equation theory deal with 100×100 sets as a commonplace and $10^3 \times 10^3$ and upwards as a *desideratum*.

Classical methods of solution by elimination or determinants have, in some hand computing at least, given way to approximate iterative techniques.

The subject, although of great antiquity, is still giving rise to a spate of research papers and there appears little agreement as to a best method of approach. The reader who wishes to pursue the subject can refer to a recent survey and bibliography⁽¹⁾ in which over 450 references to the literature are given.

7.2 DEFINITIONS

The set of simultaneous equations:

$$\left. \begin{aligned} (1) & a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n = b_1 \\ (2) & a_{21}x_1 + a_{22}x_2 + a_{23}x_3 + \dots + a_{2n}x_n = b_2 \\ (3) & a_{31}x_1 + \dots \\ & \dots \\ (r) & a_{r1}x_1 + a_{r2}x_2 + \dots + a_{rn}x_n = b_r \\ & \dots \\ (n) & a_{n1}x_1 + a_{n2}x_2 + \dots + a_{nn}x_n = b_n \end{aligned} \right\} \dots (7.2.1)$$

is taken as the basis. It may be written, in matrix form:

$$A \cdot \mathbf{x} = \mathbf{b} \dots (7.2.2)$$

DEFINITIONS

where:

$$A = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ a_{31} & \dots & \dots & a_{3n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \dots \\ x_n \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ \dots \\ b_n \end{bmatrix} \dots (7.2.3)$$

Associated with A , are the *determinant* $|A|$, the *transpose* A' [if $A = (a_{ij})$ then $A' = (a_{ji})$] and the *inverse* A^{-1} such that:

$$A^{-1} \cdot A = A \cdot A^{-1} = I \dots (7.2.4)$$

where I is the unit matrix:

$$I = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 \\ 0 & 0 & 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & \dots & 1 \end{bmatrix} \dots (7.2.5)$$

Formally, we have from equations 7.2.2 and 7.2.4:

$$A^{-1} \cdot A \cdot \mathbf{x} = A^{-1} \mathbf{b}$$

but

$$A^{-1} \cdot A \cdot \mathbf{x} = I \cdot \mathbf{x} = \mathbf{x}$$

whence

$$\mathbf{x} = A^{-1} \mathbf{b} \dots (7.2.6)$$

It is shown, in works on matrix theory⁽²⁾, that:

$$A^{-1} = \begin{bmatrix} A_{11}/|A|, & A_{21}/|A| & \dots & A_{n1}/|A| \\ A_{12}/|A|, & A_{22}/|A| & \dots & A_{n2}/|A| \\ \dots & \dots & \dots & \dots \\ A_{1n}/|A|, & A_{2n}/|A| & \dots & A_{nn}/|A| \end{bmatrix} \dots (7.2.7)$$

where A_{ij} is the co-factor of a_{ij} in $|A|$.

In the above, and in the following sections where the condition is necessary, it is assumed that $|A| \neq 0$. When this condition is satisfied, the matrix A is said to be *non-singular*.

When $\mathbf{b} = 0$ (i.e. b_r , ($r = 1 \dots n$) = 0) the equations 7.2.1 are said to be *homogeneous*; in this case the necessary and sufficient condition for non-zero solutions is $|A| = 0$.

Associated with A is the characteristic equation:

$$A \cdot x = \lambda \cdot x$$

or:

$$(A - \lambda I)x = 0 \quad \dots (7.2.8)$$

where the λ 's are scalar quantities.

Solutions of

$$|A - \lambda I| = \begin{vmatrix} a_{11} - \lambda & a_{12} & a_{13} & \dots & a_{1n} \\ a_{21} & a_{22} - \lambda & a_{23} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & \dots & a_{nn} - \lambda \end{vmatrix} = 0$$

are known as the *latent roots* of A (sometimes called *eigenwerte*, *eigenvalues*, or *characteristic numbers*). Clearly both A and A' have the same latent roots.

x and b in equations 7.2.2, 7.2.3 are seen to be composed of n components, $x_1 \dots x_n$; $b_1 \dots b_n$, and for this reason are often called *vectors*. Associated with each latent root, λ_i , say, is a *latent vector* φ_i which satisfies:

$$(A - \lambda_i I)\varphi_i = 0 \quad (i = 1 \dots n).$$

It is easily verified that $(Ax)' = x'A'$ whence, if:

$$\begin{aligned} A\varphi_i &= \lambda_i \varphi_i \\ (A\varphi_i)' &= \varphi_i' A' = \lambda_i \varphi_i' \end{aligned} \quad \dots (7.2.9)$$

Now let the latent vectors of the transpose of A , A' , be called ψ_j then, by definition:

$$A'\psi_j = \lambda_j \psi_j \quad (j = 1 \dots n)$$

whence, pre-multiplying by φ_i'

$$\varphi_i' A' \psi_j = \lambda_j \varphi_i' \psi_j,$$

similarly, post-multiplying 7.2.9 by ψ_j , we have:

$$\varphi_i' A' \psi_j = \lambda_i \varphi_i' \psi_j$$

or, subtracting the last two relationships,

$$0 = (\lambda_j - \lambda_i) \varphi_i' \psi_j.$$

It follows that, if $\lambda_j \neq \lambda_i$, $\varphi_i' \psi_j = 0$. The vectors φ_i, ψ_i are referred to as biorthogonal.

When A is symmetric, so that $A' \equiv A$, it is evident that $\varphi_i \equiv \psi_i$ ($i = 1 \dots n$). In this case $\varphi_i' \cdot \varphi_j = 0$ ($i \neq j$) and the latent vectors are orthogonal. It should be noticed that each of the latent vectors is arbitrary to the extent of a scalar multiplying constant. If these constants are so chosen that $\varphi_i' \varphi_i = 1$ ($i = 1 \dots n$), the latent vectors are said to be normalized and constitute an ortho-normal set.

Finally, with each matrix $A = (a_{ij})$ is associated a *quadratic form*

$$\sum_{i=1}^n \sum_{j=1}^n a_{ij} x_i x_j \quad (a_{ij} = a_{ji}) \quad \dots (7.2.10)$$

which is said to be *positive definite* if it is positive for every set of real (x_i) except $(x_1 = x_2 = \dots = x_n = 0)$. (Negative definiteness is defined in a similar manner).

A set of necessary and sufficient conditions for positive-definiteness is:

$$a_{11} > 0, \begin{vmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{vmatrix} > 0, \dots, \begin{vmatrix} a_{11} & \dots & a_{1n} \\ \dots & \dots & \dots \\ a_{n1} & \dots & a_{nn} \end{vmatrix} > 0. \quad \dots (7.2.11)$$

7.3 EXACT SOLUTION

The exact solution of the equations 7.2.1 is easily written down; it is:

$$x_r = \frac{\begin{vmatrix} a_{11}, a_{12} & \dots & a_{1,r-1}, b_1, a_{1,r+1} & \dots & a_{1n} \\ a_{21}, a_{22} & \dots & a_{2,r-1}, b_2, a_{2,r+1} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ a_{n1}, a_{n2} & \dots & a_{n,r-1}, b_n, a_{n,r+1} & \dots & a_{nn} \end{vmatrix}}{\begin{vmatrix} a_{11}, a_{12} & \dots & a_{1n} \\ a_{21}, a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots \\ a_{n1}, a_{n2} & \dots & a_{nn} \end{vmatrix}} \quad \dots (7.3.1)$$

Although this is formally the complete answer, it is completely useless for practical computation. The reason for this lies in the fact that, to evaluate a determinant of order n from its algebraic definition requires $(n!)(n-1)$ multiplications; it follows that the evaluation of all x_r in equation 7.3.1 would need $(n+1) \cdot (n!)(n-1)$ multiplications and n divisions. Other methods of evaluating the determinants exist, for example that of Doolittle, but these require roughly $n^3/3$ multiplications per determinant, so that the whole solution still requires about $n^4/3$ multiplications. (Notice that in all of these estimates we assume that division and multiplication are roughly equivalent in computational labour, and that n is large compared with unity.)

A more economical method is the following: consider the 1st and r th members of the set of equations 7.2.1. Multiply the former by

a_{r1} and the latter by a_{11} , and then subtract $(1) \times a_{r1}$ from $(r) \times a_{11}$. The result is:

$$(a_{11}a_{r2} - a_{r1}a_{12})x_2 + (a_{11}a_{r3} - a_{r1}a_{13})x_3 \dots$$

or

$${}_1l_{r2}x_2 + {}_1l_{r3}x_3 \dots {}_1l_{rn}x_n = {}_1b_r, \text{ say}$$

Repeating this operation for all r in the range $(1 < r \leq n)$, a set of $(n-1)$ equations in the $(n-1)$ unknowns $x_2 \dots x_n$, are obtained and a total of $2n(n-1)$ multiplications is required. (We do *not* need to form $a_{11}a_{r1}$ since this coefficient is eliminated). This sequence of operations is now repeated successively until $x_2 \dots x_{n-1}$ have each been eliminated; the result is a series of equations:

$$\left. \begin{array}{l} (1) \quad a_{11}x_1 + a_{12}x_2 + a_{13}x_3 + \dots + a_{1n}x_n = b_1 \\ (2) \quad \quad \quad 1^l_{22}x_2 + 1^l_{23}x_3 + \dots + 1^l_{2n}x_n = 1^l b_2 \\ (3) \quad \quad \quad \quad \quad 2^l_{33}x_3 + 2^l_{34}x_4 + \dots + 2^l_{3n}x_n = 2^l b_3 \\ \cdot \\ \cdot \\ \cdot \\ (n) \quad \quad \quad \quad \quad \quad \quad \quad (n-1)l_{nn}x_n = (n-1) b_n \end{array} \right\} \dots (7.3.2)$$

The total number of multiplications involved, up to this point, is

$$2[(n)(n-1) + (n-1)(n-2) + \dots + 2 \cdot 1] = \frac{2}{3}n(n^2 - 1)$$

Starting with the n th member of equation 7.3.2 we can find $x_n, x_{n-1}, x_{n-2} \dots x_1$ successively by substitution of the values of x_r already found in the next equation [*i.e.* x_n in $(n-1)$ to get x_{n-1} ; x_n, x_{n-1} in $(n-2)$ to get x_{n-2} etc.]. It will be seen that this process involves $\frac{1}{2}n(n-1)$ multiplications and n divisions, whence the total number of 'equivalent multiplications' for the whole solution is:

$$\frac{2}{3}n(n^2 - 1) + \frac{1}{3}n(n - 1) + n$$

which is of the order $\frac{2}{3}n^3$ for large values of n .

The merit of this process lies in the fact that, if the initial coefficients a_{ij} are small whole numbers, no round-off may be necessary during the initial elimination; thus if integer solutions exist they will be obtained correctly. (This may be important if the equations are 'ill conditioned,' *vide infra*.)

When the initial coefficients are *not* integers a more efficient process is as follows,

- (1) Examine the matrix $A = [a_{ij}]$ and find the largest a_{ij} , call this a_{IJ} .
- (2) Multiply the I th equation successively by all a_{iJ}/a_{IJ} ($i = 1, \dots, n \neq I$) and subtract from each of the other equations.

The result of this operation will be a set of $(n - 1)$ equations from which x_r has been eliminated, just as occurred in the previous method. The same discrimination-elimination process is applied to the new set of equations and continued until a triangular set of equations of the same type as 7.3.2 is obtained. Back substitution then yields the x_r .

The reason for the choice of the *largest* a_{ij} for the divisor at each stage is that the round-off errors, generated at each formation of a_{IJ}/a_{IJ} , are thereby reduced in size at the next stage. The method is sometimes described as 'pivotal condensation,' the element a_{IJ} being the pivot.

It is easy to see that $n(n-1)$ multiplications and $(n-1)$ divisions have to be made at the first elimination. This gives (n^2-1) 'equivalent multiplications,' compared to $2n(n-1)$ in the previous process. Thus, for the whole solution, a number of multiplications of order $\frac{1}{2}n^3$ is required.

For hand computation, pivotal condensation is an excellent method since the largest coefficient a_{IJ} can usually be seen at a glance. When an automatic computing machine is in use, however, this discrimination may be troublesome and, for this reason, it has been our practice to use the former method whenever possible.

A third method of solution is usually known as the Choleski process; it depends upon the reduction of the matrix A to the product of a lower triangular matrix L , with unit diagonal coefficients, and an upper triangular matrix U .

Thus: $Ax = L.U.x = b$

whence if

$$L\xi = b$$

$$Ux = \xi$$

and the solution is reduced to a pair of back substitutions. The coefficients l_{ij} and u_{ij} are readily obtained by forming the product L, U and equating coefficients to those of A in the sequence

$$(a_{11} \dots a_{1n})(a_{21} \dots a_{2n}) \dots (a_{r1} \dots a_{rn}) \dots (a_{n1} \dots a_{nn}).$$

It is stated that an advantage of the Choleski method lies in the fact that few intermediate results have to be written down; it is our

experience, however, that the method is too complicated for 'occasional' use, and is not suitable for use on an automatic computer.

7.4 THE INVERSION OF MATRICES

When it is desired to solve the equations 7.2.1 for a fixed set of a_{ij} , but for a number of vectors b , it is best to evaluate the inverse matrix A^{-1} and to use equation 7.2.6 to obtain the individual solutions.

The formal inverse (7.2.7) is useless, for the same reasons which cause the rejection of the determinant solution of the original equations. The elimination method may, however, be extended to give the inverse of a matrix. Thus let:

$$A^{-1} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \quad \dots (7.4.1)$$

Then:

$$\begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2n} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \\ \dots \\ b_n \end{bmatrix} = \begin{bmatrix} a_{11}b_1 + a_{12}b_2 + \dots + a_{1n}b_n \\ a_{21}b_1 + a_{22}b_2 + \dots + a_{2n}b_n \\ \dots \\ a_{n1}b_1 + a_{n2}b_2 + \dots + a_{nn}b_n \end{bmatrix} \quad \dots (7.4.2)$$

Whence, if we solve the original set of equations for the vectors:

$$b = \begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \dots \\ 0 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \\ \dots \\ 0 \end{bmatrix} \dots \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \dots \\ 1 \end{bmatrix} \quad \dots (7.4.3)$$

the solutions, considered as column vectors, are the columns of the inverse matrix.

We note that by using pivotal condensation to solve the sets of equations and arranging the work so that only one triangulation of A is performed $\frac{2}{3}n^3$ multiplications are required.

The Choleski method may also be applied to matrix inversion, for let

$$A = L.U$$

where L and U are triangular matrices of the type defined in section 7.3. Then:

$$A.A^{-1} = I$$

$$\text{whence: } L.U.A^{-1} = I$$

$$L^{-1}LUA^{-1} = UA^{-1} = L^{-1}I = L^{-1}$$

$$U^{-1}UA^{-1} = A^{-1} = U^{-1}L^{-1}.$$

Now it has been seen in section 7.3 that U and L can be determined from A ; assume that:

$$L^{-1} = \begin{bmatrix} \lambda_{11} & \lambda_{12} & \dots & \lambda_{1n} \\ \lambda_{21} & \lambda_{22} & \dots & \lambda_{2n} \\ \dots & \dots & \dots & \dots \\ \lambda_{n1} & \lambda_{n2} & \dots & \lambda_{nn} \end{bmatrix}, \quad L = \begin{bmatrix} 1 & 0 & 0 & \dots & 0 \\ l_{21} & 1 & 0 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots \\ l_{n1} & l_{n2} & \dots & l_{n,n-1} & 1 \end{bmatrix}.$$

Then, since $L.L^{-1} = I$, we have, equating coefficients:

$$\lambda_{11} = 1, \lambda_{12} = \lambda_{13} = \dots = \lambda_{1n} = 0$$

$$\lambda_{11}l_{21} + \lambda_{21} = 0, \quad \lambda_{12}l_{21} + \lambda_{22} = 1,$$

$$\lambda_{13}l_{21} + \lambda_{23} = 0 \dots \lambda_{1j}l_{21} + \lambda_{2j} = 0 (j > 2) \text{ etc.}$$

which will enable the λ_{ij} to be determined. In a like manner the components Y_{ij} of the inverse U^{-1} may be determined. Thus, multiplying $U^{-1}L^{-1}$, the inverse A^{-1} is obtained.

Two other methods of finding an inverse are worthy of mention, not so much for their efficiency as far as numbers of multiplications are concerned, but because they are very well adapted to punched card and automatic computing machinery.

The first method is a consequence of the Cayley-Hamilton theorem ('A matrix satisfies its characteristic equation').

$$\text{Let } \lambda^n + a_1\lambda^{n-1} + a_2\lambda^{n-2} \dots + a_n = 0 \quad \dots (7.4.4)$$

be the characteristic equation of a matrix A . Then:

$$A^n + a_1A^{n-1} + a_2A^{n-2} + \dots + a_nI = 0$$

$$\text{whence: } I = -\frac{1}{a_n} (A^n + a_1A^{n-1} + a_2A^{n-2} \dots + a_{n-1}A)$$

and:

$$A^{-1}I = A^{-1} = -\frac{1}{a_n} (A^{n-1} + a_1A^{n-2} + a_2A^{n-3} + \dots + a_{n-1}I). \quad \dots (7.4.5)$$

Now it may be shown that:

$$\left. \begin{aligned} a_1 &= -s_1 \\ a_2 &= -\frac{1}{2}(a_1s_1 + s_2) \\ a_3 &= -\frac{1}{3}(a_2s_1 + a_1s_2 + s_3) \\ &\dots\dots\dots \\ a_n &= -\frac{1}{n}(a_{n-1}s_1 + a_{n-2}s_2 \dots + s_n) \end{aligned} \right\} \dots\dots(7.4.6)$$

where:

$$s_r = \text{tr}(A^r) \dots\dots(7.4.7)$$

tr being the *trace* of the matrix (i.e. $\sum_{i=1}^n a_{ii}$ the sum of the elements on the principal diagonal).

Since a single matrix multiplication requires n^3 scalar multiplications the process of inversion embodied in equations 7.4.5, 7.4.6 and 7.4.7 will require of the order of n^4 multiplications.

The second inversion technique is an iterative one based upon the well-known result:

$$x_{n+1} = x_n(2 - ax_n) \dots\dots(7.4.8)$$

for which

$$\text{Lt}_{n \rightarrow \infty} x_n = 1/a$$

Under suitable conditions⁽³⁾ it may be shown that a similar iteration:

$$X_{n+1} = X_n(2I - AX_n) \dots\dots(7.4.9)$$

in which X_n and A are matrices and I is the unit matrix, will converge to A^{-1} . A practical experiment⁽⁴⁾, in which this method was applied to the inversion of a 16×16 matrix, showed that 10 iterations were required to produce 4 decimal place accuracy and 20 iterations to produce 6 decimal place accuracy. Actually, the rate of convergence depends upon the goodness of the original approximation, X_0 , and upon the precise form of A . It has been suggested that $X_0 = I$ should be taken as a starting value.

7.5 RESIDUALS AND 'CONDITION'

With the exception of the last method described, the procedures so far outlined for solving a set of linear simultaneous equations will give an *exact* solution if the arithmetic operations can be performed with complete accuracy. In practice only an approximation will, in

general, be produced, since the arithmetic operations \times and \div will be subject to round-off errors. The iterative methods which will now be described make no pretence at complete accuracy, but merely reduce the errors in an existing approximation. This leads naturally to a consideration of the way in which the accuracy of a given approximation may be measured.

Let us suppose that an approximation $\xi (= \xi_1, \xi_2, \dots, \xi_n)$ is given to the solution of the set of equations 7.2.1. If we substitute this in the given equations we obtain:

$$\left. \begin{aligned} a_{11}\xi_1 + a_{12}\xi_2 + \dots + a_{1n}\xi_n - b_1 &= r_1 \\ a_{21}\xi_1 + a_{22}\xi_2 + \dots + a_{2n}\xi_n - b_2 &= r_2 \\ &\dots\dots\dots \\ a_{n1}\xi_1 + a_{n2}\xi_2 + \dots + a_{nn}\xi_n - b_n &= r_n \end{aligned} \right\} \dots\dots(7.5.1)$$

The quantities r_i define a vector $\mathbf{r} (= r_1, r_2, \dots, r_n)$ which is a measure of the inaccuracy of the approximation. To obtain a single number which will express the inaccuracy we may consider the *length* of \mathbf{r} , or more usefully its square, R^2 . This is defined as the scalar product of \mathbf{r} and itself:

$$R^2 = (\mathbf{r}, \mathbf{r}) = \sum_{i=1}^n r_i^2 = (A\xi - \mathbf{b}) \cdot (A\xi - \mathbf{b}) \dots\dots(7.5.2)$$

and the quantities r_i are often called 'residuals.'

By means of equation 7.5.1 we may reduce 7.5.2 to a quadratic form in the ξ_i :

$$R^2 = \sum_{i=1}^n (a_{i1}\xi_1 + a_{i2}\xi_2 \dots + a_{in}\xi_n - b_i)^2 \dots\dots(7.5.3)$$

which is clearly positive-definite.

An alternative measure of accuracy is the quadratic (previously given in 7.2.10):

$$S = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_{ij}\xi_i\xi_j - \sum_{i=1}^n b_i\xi_i \quad (a_{ij} = a_{ji}) \dots\dots(7.5.4)$$

This has the merit of bearing a simple relationship to the original matrix A , but the disadvantage that A must be positive definite, so that arguments based upon S may fail if this condition is not satisfied. We note, here, that equation 7.5.4 may be written, vectorially:*

* Notice that in equations 7.5.2, 7.5.5, A has the character of a 'Tensor,' or in more modern usage, a 'Dyadic.'

$$S = \frac{1}{2} \xi \cdot A \xi - b \cdot \xi \quad \dots (7.5.5)$$

$$= \frac{1}{2} \xi \cdot (r + b) - b \cdot \xi$$

$$= \frac{1}{2} \xi \cdot (r - b)$$

$$= \frac{1}{2} (r - b) \cdot A^{-1} (r + b) \quad \dots (7.5.6)$$

It is clear from equation 7.5.4 that S has a *minimum* when $r = 0$, and the same is true of R^2 .

The question now arises, how far may an approximate solution ξ deviate from the true solution x for a particular value of R^2 or S . A measure of this deviation is the square of the length of the difference vector $(\xi - x)$, that is

$$(\xi - x) \cdot (\xi - x) = (A^{-1}r) \cdot (A^{-1}r).$$

This is not easily correlated with either R^2 or S but it shows that if A^{-1} has any large component, large differences between true and approximate solutions may be accompanied by small residuals r_i and values of R^2 and S .

Another approach is to consider the hyper-ellipsoids defined by equations 7.5.3 and 7.5.4, for *fixed* values of R^2 and S . Considering the latter (S), we first transform to axes through the centre $(x_1, x_2 \dots x_n)$. Putting

$$\xi_r = x_r + \epsilon_r$$

equation 7.5.4 becomes:

$$S + \frac{1}{2} \sum_{i=1}^n b_i x_i = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n a_{ij} \epsilon_i \epsilon_j. \quad \dots (7.5.7)$$

Now it is well known ⁽⁵⁾ that the quadratic form:

$$\sum_{ij} a_{ij} \epsilon_i \epsilon_j$$

can, by means of a real orthogonal transformation of unit modulus, be reduced to:

$$\sum_{i=1}^n \lambda_i E_i^2 \quad \dots (7.5.8)$$

where the λ_i are the roots of the characteristic equation 7.2.9.

Now

$$S + \frac{1}{2} \sum_{i=1}^n b_i x_i = \sum_{i=1}^n \lambda_i E_i^2$$

represents a hyper-ellipsoid having axes proportional to $1/\sqrt{\lambda_i}$

so that, if any of the λ_i are very small, large values of E_i can be

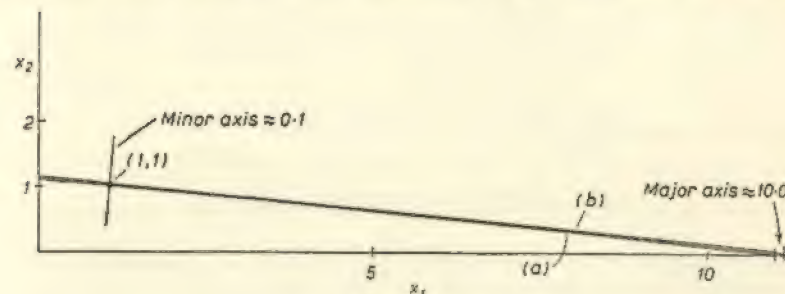


Figure 7.5.1

accompanied by small residuals. This is shown, for a two-dimensional case, in Figure 7.5.1.

The equations considered are:

$$x_1 + 10x_2 = 11 \quad \dots (a)$$

$$10x_1 + 101x_2 = 111 \quad \dots (b)$$

for which the solution is obviously $x_1 = x_2 = 1$;

Figure 7.5.1 shows the lie of the curve for $S = 60.5$, the actual curve being too close to its major axis to be visible on this scale. It will be seen that the axes are approximately $\frac{1}{10}$ and 10 in accord with the ratio $\sqrt{\lambda_1/\lambda_2}$ of the roots of:

$$\begin{vmatrix} 1 - \lambda & 10 \\ 10 & 101 - \lambda \end{vmatrix} = 0$$

or

$$\lambda^2 - 102\lambda + 1 = 0$$

giving

$$\lambda_1 = 101.9902, \quad \lambda_2 = 0.0098.$$

Also plotted are the lines represented by the equations (a) and (b), and it will be noticed that they are nearly parallel, so that the intersection (1, 1) is very sensitive to the coefficients. The relationship between residuals and parallelism is the basis of the definition of R^2 in equation 7.5.2; thus if p_i is the length of the perpendicular from $(\xi_1, \xi_2 \dots \xi_n)$ on to the hyperplane:

$$a_{i1}x_1 + a_{i2}x_2 + a_{i3}x_3 \dots a_{in}x_n - b_i = 0$$

we have:

$$p_i = \left(\sum_{j=1}^n a_{ij} \xi_{j1} - b_i \right) / \sqrt{\sum_{j=1}^n a_{ij}^2} \quad \dots (7.5.9)$$

$$= r_i / \sqrt{\sum_{j=1}^n a_{ij}^2}$$

If now we normalize the rows of the matrix A

$$\left(\text{i.e., divide each row by } \sqrt{\sum_{j=1}^n a_{ij}^2} \right)$$

we shall not alter the solution of the equations. Assume that this has been done, so that equation 7.5.9 becomes:

$$p_i = r_i. \quad \dots (7.5.10)$$

Thus equation 7.5.2 may be written:

$$R^2 = \sum_{i=1}^n p_i^2$$

and we see that the surfaces of $R^2 = \text{const.}$ are simply those such that the sums of the squares of the perpendiculars from any point $(\xi_1, \xi_2, \dots, \xi_n)$ on the surface on to the hyperplanes

$$a_{i1}x_1 + a_{i2}x_2 + \dots + a_{in}x_n - b_i = 0 \quad (i = 1 \dots n)$$

is constant.

It will be evident, from *Figure 7.5.1*, that in the example quoted, the $R^2 = \text{const.}$ 'surfaces' are similar in shape to those defined by $S = \text{const.}$

The equations (a) and (b) provide an interesting study of the false conclusions which may be drawn from small values of R^2 . We first normalize by forming (a)/ $\sqrt{1^2 + 10^2}$, (b)/ $\sqrt{10^2 + 101^2}$ and thus obtain:

$$\frac{x_1}{\sqrt{101}} + \frac{10x_2}{\sqrt{101}} = \frac{11}{\sqrt{101}} \dots \dots \dots (a')$$

$$\frac{10x_1}{\sqrt{10301}} + \frac{101x_2}{\sqrt{10301}} = \frac{111}{\sqrt{10301}} \dots \dots (b').$$

The radicals have been purposely retained thus far. If we calculate the components of the residuals for the approximate solution $x_1 = 1.001$, $x_2 = 1.01$ we find: from (a') $r_1 \approx .01$, from (b') $r_2 \approx .01$ whence $R^2 \approx .0002$.

Next try the values $x_1 = 11.1$, $x_2 = 0$, we obtain: $r_1 \approx .01$, $r_2 = 0$ whence, this time $R^2 \approx .0001$. Most people would conclude from this

that $(11.1, 0)$ was 'nearer' to the true solution than $(1.001, 1.01)$ which is manifestly false.

If we had evaluated the radicals and expressed the resulting coefficients in decimal form, the equations (a') and (b') would have appeared as:

$$.09950x_1 + .99504x_2 = 1.09454 \quad (a'')$$

$$.09853x_1 + .99513x_2 = 1.09366 \quad (b'')$$

and a change of approximately 1 per cent in the coefficient of x_1 in (a'') makes the solution $(11.1, 0)$ exact.

Thus, the solution of equations of this type is very sensitive to the values of the coefficients, and the term 'ill-conditioned' is applied to them.

Various measures for 'ill-condition' suggest themselves from the preceding discussion. Perhaps the most obvious is the ratio $\lambda_{\max}/\lambda_{\min}$, which gives a measure of the ratio of greatest to least axes of the hyper-ellipsoid defined by equation 7.5.4. Unfortunately the calculation of λ_{\max} and λ_{\min} is an operation of at least the complexity of the solution of the original equations, so that this criterion is of little practical use. Ill condition is accompanied by approximate parallelism of some of the hyperplanes defined by the equations. This can be recognized in simple cases, such as that obtaining in our equations (a) and (b), by observing the ratios of corresponding coefficients, but this, too, is not generally possible.

Yet another measure which has been suggested is the value of the determinant $|A|$, 'small' values being accompanied by ill-condition. In this form the test is without value, since the set of equations can be multiplied by any constant without altering the solution. If the constant is M , $|A|$ is multiplied by M^n and so can be made as large as desired; we note that if $|A| = 0$ the equations have no unique solution. A better version of this test is first to 'normalize' the equations by division by $\sqrt{\sum_{j=1}^n a_{ij}^2}$ for the i th equation, and to regard smallness compared to ± 1 as an indication of ill-condition. On this basis the 'ideal' equations:

$$a_{11}x_1 = b_1$$

$$a_{22}x_2 = b_2$$

and

$$a_{11}x_1 + a_{12}x_2 = b_1$$

$$a_{21}x_1 - a_{22}x_2 = b_2$$

with

$$a_{11} = a_{12}, \quad a_{21} = a_{22}$$

have $|A| = 1$ $|A| = -1$ respectively, when normalized. On the other hand, our equations (a) (b), in the normal form (a'') (b''), have $|A| \approx .001$ which gives a clear indication of their ill-condition. Yet another example is provided by the much quoted⁽⁶⁾ set of ill-conditioned equations:

$$\begin{aligned} 5x_1 + 7x_2 + 6x_3 + 5x_4 &= 23 \\ 7x_1 + 10x_2 + 8x_3 + 7x_4 &= 32 \\ 6x_1 + 8x_2 + 10x_3 + 9x_4 &= 33 \\ 5x_1 + 7x_2 + 9x_3 + 10x_4 &= 31 \end{aligned}$$

whose true solution is evidently (1, 1, 1, 1), but for which the values (+14.6, -7.2, -2.5, +3.1) give components (+.1, -.1, -.1 + .1) for the residue vector r , and consequently $R^2 = .04$. The values (+2.36, +0.18, +0.65, +1.21) give residue components (+.01, -.01, -.01, +.01) and consequently $R^2 = .0004$. The value of $|A|$ for this set of equations as it stands is 1, but if we normalize by division of the respective equations by

$$\begin{aligned} \sqrt{5^2 + 7^2 + 6^2 + 5^2}, \quad & \sqrt{7^2 + 10^2 + 8^2 + 7^2}, \\ \sqrt{6^2 + 8^2 + 10^2 + 9^2} \quad \text{and} \quad & \sqrt{5^2 + 7^2 + 9^2 + 10^2}, \end{aligned}$$

we get $|A|_{\text{norm.}} = .000,019,9$ which, since it is small compared with unity, indicates clearly the ill condition of the equations.

7.6 DETERMINATION OF LATENT ROOTS AND OF CHARACTERISTIC VECTORS

It is possible to determine the latent roots of a matrix directly from the definition 7.2.8, and then to solve the resulting sets of simultaneous equations (one set for each λ_i) to determine the characteristic vectors. This method, whilst sometimes appropriate for small numbers of equations and unknowns, say less than 5, is not at all feasible when large numbers of variables are involved. This point is therefore an appropriate one at which to consider methods of iteration and successive approximation which will, at the same time, be useful for the solution of simultaneous equations themselves.

When no information is available regarding either latent roots or characteristic vectors, it is usual to start by determining the largest λ_i and its associated vector.

We shall adopt the usual procedure and denote the latent roots by $\lambda_1, \lambda_2 \dots \lambda_n$, where $|\lambda_1| \geq |\lambda_2| \geq |\lambda_3| \dots \geq |\lambda_n|$, the associated characteristic vectors are $\varphi_1, \varphi_2 \dots \varphi_n$ and we recall the definition:

$$(A - \lambda_i I)\varphi_i = 0 \quad (i = 1 \dots n) \quad \dots (7.6.1)$$

Consider the effect of applying the tensor A to any vector x , that is, of forming Ax . The tensor operator can be considered as a set of expansions, or contractions, in an n dimensional hyperspace. These expansions have *directions* given by the principal axes of the quadratic form $\sum_{i,j} a_{ij}x_i x_j$, that is, by the directions of the characteristic vectors; their magnitudes are such that a unit vector in the direction of the i th principal axis has its length changed to λ_i . In two dimensions the effect is shown in Figure 7.6.1.

This was constructed for latent roots $\lambda_1 = 2, \lambda_2 = (\frac{1}{2})$ and characteristic vectors φ_1, φ_2 in the directions shown. Points (1), (2), (3), (4)

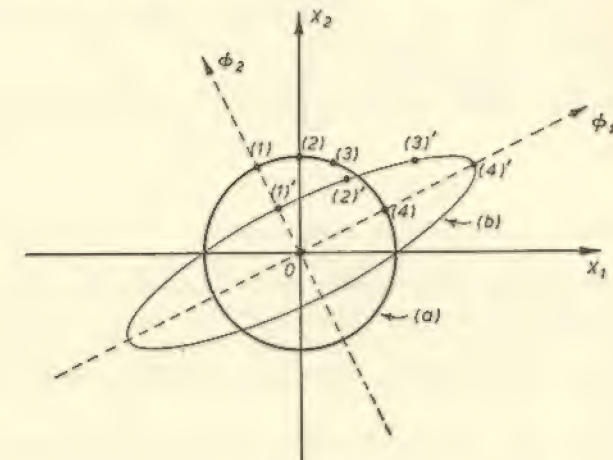


Figure 7.6.1

on the unit circle (a) are transformed into (1)', (2)', (3)', (4)' by the tensor.

Next consider the effect of forming Ax , normalizing to form $(Ax)_n$, and then reapplying A .

This is seen, from Figure 7.6.2, to result in a vector in which the component in direction φ_1 is increased by a factor λ_1 , and that in direction φ_2 by λ_2 . Thus the transformed vector is more nearly in the direction φ_1 than was originally the case. The sequence followed by successive applications of the process to a unit vector, initially (1), is shown. $A(1) = (1)'$. This is normalized to give $(1)_n$, $A(1)_n$ gives $(2)'$ and so on. It is clear that a few applications of the process will lead to a vector which lies, sensibly, in direction φ_1 . The figure also shows that, if an initial vector (0) had been chosen *orthogonal* to φ_1 ,

the operation $A(0)$ would not lead to any convergence to φ_1 , but merely to oscillation along φ_2 . The whole argument can be generalized immediately to an n dimensional hyperspace.

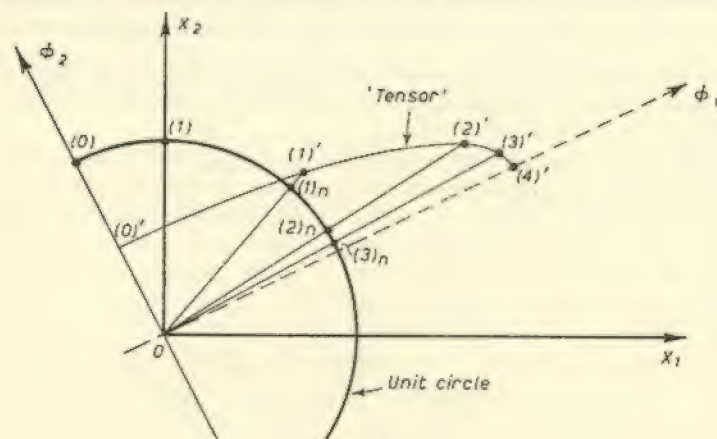


Figure 7.6.2

An alternative demonstration of the process just described, if A is symmetrical (i.e. $a_{ij} = a_{ji}$), is as follows. The characteristic vectors are orthogonal, that is $(\varphi_i, \varphi_j) = 0$ ($i \neq j$); we may therefore express an arbitrary vector \mathbf{x} as a linear combination of the φ_i :

$$\mathbf{x} = a_1\varphi_1 + a_2\varphi_2 + \dots + a_n\varphi_n.$$

Now apply the operator A and use equation 7.6.1:

$$\begin{aligned} A\mathbf{x} &= a_1\lambda_1\varphi_1 + a_2\lambda_2\varphi_2 + \dots + a_n\lambda_n\varphi_n \\ &= \lambda_1 \left(a_1\varphi_1 + a_2\frac{\lambda_2}{\lambda_1}\varphi_2 + \dots + a_n\frac{\lambda_n}{\lambda_1}\varphi_n \right) \end{aligned}$$

so that, since $\lambda_1 > \lambda_n$ ($n > 1$) the result of the operation A is to decrease the proportion of vectors, other than φ_1 , in the system. Thus the process converges, eventually, to $\lambda_1\varphi_1$.

In practice it is usually simpler to take the vector \mathbf{x} as having a maximum component unity instead of normalizing properly. This avoids the evaluation of $(\sum a_n^2)^{-\frac{1}{2}}$ at each stage, and does not otherwise affect the process.

When multiple roots occur, so that $\lambda_1 = \lambda_2 = \dots = \lambda_r$ (say), there is no unique value of φ_1 and, under these circumstances the method breaks down. If, in practical application, it is found that

convergence is either slow or absent, it is well, before abandoning the process, to try an initial vector which is orthogonal to that originally taken, in order to eliminate the possibility that the first choice was at right angles to φ_1 .

The figure taken as an illustration in the preceding paragraphs was constructed on the basis of a positive definite quadratic form $\sum a_{ij}x_ix_j$, for which all λ_i are greater than zero. The same argument is applicable when some λ_i are negative, with the exception that the end points (1)', (2)' etc. will move from side to side of φ_2 (assuming $|\lambda_1| > |\lambda_2|$ and λ_1 negative). The normalizing process now consists not only of reducing the amplitude of $A\mathbf{x}$ to unity, but also in restoring it to its original side of φ_1 .

For situations in which convergence is slow, various modifications of the above process are possible. In the first, instead of successively multiplying the arbitrary starting vector by A , some high power of A is first formed by the sequence:

$$A, A^2, A^4, A^8, A^{16} \dots \text{etc.}$$

Evidently, to form A^{2^n} requires n matrix multiplications, so that quite a small number of multiplications produce a high power of A . The iterative procedure then uses A^{2^n} instead of the simple A . It should be noted, however, that each matrix multiplication needs about $\frac{1}{2}n^3$ arithmetical multiplications so that detailed consideration needs to be given to work estimates before one method can be described as better than the other.

Suppose that the ordinary method of iteration with A requires k steps and that A is an $n \times n$ matrix. The number of arithmetical multiplications required is thus kn^2 . If now, the power of A needed to produce the same result in one step is 2^m we have that $2^m = k$ and that the number of arithmetical multiplications needed to generate A^{2^m} is $\frac{1}{2}mn^3$.

It follows that, for the power method to be efficient,

$$\frac{1}{2}mn^3 < 2^mn^2$$

or

$$m < 2^{m+1}/n$$

Some typical values are shown below.

Order of matrix, n	4	8	16
m	≥ 3	≥ 5	≥ 6
$k = 2^m =$ number of multiplications required for simple iterative method	≥ 8	≥ 32	≥ 64

The values of k may appear too large, but when it is noticed that if the dominant and subdominant roots have a ratio A , say 0.9, to attain an accuracy of 10^{-6} in φ_1 will, in general, need about 130 iterations with A , the advantage of the power method becomes obvious.

Another method of improving convergence is to operate with $(A - pI)$ rather than with A . Thus

$$(A - pI)\varphi_i = (\lambda_i - p)\varphi_i$$

so that $(A - pI)$ has the same latent vectors φ_i as A but has as latent roots $(\lambda_i - p)$. Consider, for example, a matrix whose latent roots are $+3 + 2$ and $+1$. The rate of convergence of the iterative method with an arbitrary starting vector will not be better than a factor $(\frac{2}{3})$ at each stage. If, however, we iterate with $(A - I)$ the rate of convergence will be proportional to $(\frac{1}{2})$ and if with $(A - \frac{3}{2}I)$ to $(\frac{1}{3})$, an improvement by a factor of 2 over the original.

A third method of speeding up the determination of latent roots and vectors by the iterative method is to use Aitken's δ^2 process. Assume that three successive iterates u_0 , u_1 and u_2 are known and that u_2 tends to a limit in such a manner that the successive terms form an approximately geometric progression. In this event it may be assumed that:

$$u_0 = u_{\infty} + kr^n$$

$$u_1 = u_{\infty} + kr^{n+1}$$

$$u_2 = u_{\infty} + kr^{n+2}$$

$$\text{whence } \frac{u_0 - u_{\infty}}{u_1 - u_{\infty}} = \frac{u_1 - u_{\infty}}{u_2 - u_{\infty}}$$

$$\text{or } u_{\infty} = u_0 u_2 - u_1^2 / u_0 - 2u_1 + u_2$$

It will be found that computation of the numerator of this expression usually involves differencing two nearly-equal quantities and thus implies high accuracy in the generation of the products $u_0 u_2$ and u_1^2 . This difficulty is avoided by writing

$$u_{\infty} = u_2 - \frac{(u_2 - u_1)^2}{u_0 - 2u_1 + u_2}$$

As an example of the power of the Aitken process, consider the matrix:

$$A = \begin{bmatrix} -5.509\ 882 & 1.870\ 086 & 0.422\ 908 \\ 0.287\ 865 & -11.811\ 654 & 5.711\ 900 \\ 0.049\ 099 & 4.308\ 033 & -12.970\ 687 \end{bmatrix}$$

Starting with $x_1 = [1, 0, 0]'$, we find that:

x_{12}	Ax_{12}	x_{13}	Ax_{13}	x_{14}	Ax_{14}
1.000 000	-17.351 78	1.000 000	-17.378 48	1.000 000	-17.389 55
-8.113 909	141.126 75	-8.133 271	141.483 89	-8.141 326	141.632 99
7.873 325	-137.093 17	7.900 812	-137.468 26	7.910 257	-137.625 47

It is seen that convergence is relatively slow but, applying the δ^2 process to the values:

$$-17.35178, \quad -17.37848, \quad -17.38955$$

we at once find:

$$\lambda_{\max} = 17.3974$$

which is a much better approximation to the true $\lambda_{\max} = 17.3977$ than any of the original values.

When λ_1 is complex a variant of the above process is adopted. After a sufficient number of iterations it can be shown that

$$A^n x + bA^{n-1}x + cA^{n-2}x \rightarrow 0$$

where b and c are numbers such that

$$\lambda^2 + b\lambda + c = 0$$

is the quadratic equation satisfied by λ_1 and its complex conjugate.

Thus, by choosing two values of n and remembering that vector equations are also satisfied by individual components, pairs of linear simultaneous equations are obtained which give b and c , and consequently λ_1 .

An alternative, and possibly simpler, method has been suggested by workers at the National Physical Laboratory. This depends upon iteration with the operator $(A + ipI)$ where p is a suitably chosen real number.

When an approximation is known to any characteristic vector of a symmetric matrix, it is possible to obtain an improved value for the associated latent root by means of the expression:

$$L = \frac{(x \cdot Ax)}{(x \cdot x)} = \frac{\sum_{i,j} a_{ij} x_i x_j}{\sum_i x_i^2} \dots (7.6.2)$$

Assume that

$$x = \varphi_i + \sum_{j \neq i} \epsilon_j \varphi_j$$

where the φ from an orthonormal set and the ϵ are small, scalar multipliers. Then,

$$\begin{aligned} \mathbf{x}' \cdot \mathbf{A} \mathbf{x} &= (\varphi'_i + \sum_{j \neq i} \epsilon_j \varphi'_j)' A (\varphi_i + \sum_{j \neq i} \epsilon_j \varphi_j) \\ &= (\varphi'_i + \sum_{j \neq i} \epsilon_j \varphi'_j)' (\lambda_i \varphi_i + \sum_{j \neq i} \epsilon_j \lambda_j \varphi_j) \\ &= \lambda_i + \sum_{j \neq i} \epsilon_j^2 \lambda_j \end{aligned}$$

similarly $(\mathbf{x}' \cdot \mathbf{x}) = (\varphi_i + \sum_{j \neq i} \epsilon_j \varphi_j)' (\varphi_i + \sum_{j \neq i} \epsilon_j \varphi_j)$

$$= 1 + \sum_{j \neq i} \epsilon_j^2$$

whence
$$L = \lambda_i \left(1 + \sum_{j \neq i} \frac{\lambda_j}{\lambda_i} \epsilon_j^2 \right) / \left(1 + \sum_{j \neq i} \epsilon_j^2 \right)$$

$$= \lambda_i \left(1 + \sum_{j \neq i} \left(\frac{\lambda_j}{\lambda_i} - 1 \right) \epsilon_j^2 \right) + O(\epsilon^4)$$

so that the error in λ_i determined from L , is of order ϵ^2 .

We may notice, before leaving the quantity L , that if \mathbf{x} is any vector whatever, and A is a real symmetric square matrix, then:

$$\lambda_1 \geq L \geq \lambda_n. \quad \dots (7.6.3)$$

Having determined the largest latent root and its associated characteristic vector, it is possible to proceed to the determination of the next smaller root and vector. A crude method is to select any vector ${}_1\varphi_2$ orthogonal to φ_1 and then to go through the iterative process $A_1\varphi_2, A(A_1\varphi_2)_n$ etc. Since ${}_1\varphi_2$ is orthogonal to φ_1 we have seen that convergence to φ_1 will not occur, and instead the process will, in principle, tend to $\lambda_2\varphi_2$. Unfortunately, since the arithmetic processes used to form $A_1\varphi_2$ etc. are subject to round off errors, a small proportion of φ_1 will be introduced into the ${}_m\varphi_2$ vectors, and this will grow until convergence to φ_1 instead of φ_2 results.

This defect is removed in an iterative method which is based upon the Gram-Schmidt process for the construction of orthogonal vectors.

Assume that φ_1 is known, we then take any vector ${}_1\varphi_2$ and form:

$$A_1\varphi_2 - a_1\varphi_1$$

so that this vector is orthogonal to φ_1 . Now let:

$${}_2\lambda_2 \cdot {}_2\varphi_2 = A \cdot {}_1\varphi_2 - a_1\varphi_1$$

and repeat the process to obtain successively:

$$\begin{aligned} {}_3\lambda_2 \cdot {}_3\varphi_2 &= A \cdot {}_2\varphi_2 - a_2\varphi_1 \\ {}_4\lambda_2 \cdot {}_4\varphi_2 &= A \cdot {}_3\varphi_2 - a_3\varphi_1 \\ &\dots \dots \dots \\ {}_m\lambda_2 \cdot {}_m\varphi_2 &= A \cdot {}_{m-1}\varphi_2 - a_{m-1}\varphi_1. \end{aligned}$$

The multipliers ${}_2\lambda_2, {}_3\lambda_2$, etc. are inserted in this scheme to show that the vectors ${}_2\varphi_2, {}_3\varphi_2$ etc. are in standardized form. This standardization may be effected either by true normalization so that the sum of the squares of the components of the ${}_n\varphi_2$ are unity or, more easily, by scaling so that the leading component is made equal to unity at each stage.

The scalar multipliers a_1, a_2 etc. should be small after the first stage, and serve merely to keep φ_1 from contaminating the convergents to φ_2 . It is probably sufficient, in most cases, to insert an $a\varphi_1$ term, not at each stage, but once every 5 cycles. The values of a_1, a_2 etc. are determined from the orthogonality of φ_1 and each of the ${}_n\varphi_2$, thus

$${}_n\lambda_2 \varphi'_1 \cdot {}_n\varphi_2 = 0 = \varphi'_1 \cdot A \cdot {}_{n-1}\varphi_2 - a_{n-1} \varphi'_1 \cdot \varphi_1$$

whence
$$a_{n-1} = \varphi'_1 \cdot A \cdot {}_{n-1}\varphi_2 / \varphi'_1 \cdot \varphi_1$$

When λ_2 and φ_2 have been determined with sufficient accuracy, we chose a vector φ_3 and form:

$$\begin{aligned} {}_2\lambda_3 \cdot {}_2\varphi_3 &= A \cdot {}_1\varphi_3 - a_1\varphi_1 - \beta_1\varphi_2 \\ {}_3\lambda_3 \cdot {}_3\varphi_3 &= A \cdot {}_2\varphi_3 - a_2\varphi_1 - \beta_2\varphi_2 \\ &\dots \dots \dots \\ {}_m\lambda_3 \cdot {}_m\varphi_3 &= A \cdot {}_{m-1}\varphi_3 - a_{m-1}\varphi_1 - \beta_{m-1}\varphi_2. \end{aligned}$$

The values of $a_1, a_2, \beta_1, \beta_2$ etc. are again determined from the orthogonality of φ_1 and φ_2 to each of the ${}_n\varphi_3$. Thus,

$${}_n\lambda_3 \cdot \varphi'_1 \cdot {}_n\varphi_3 = 0 = \varphi'_1 A \cdot {}_{n-1}\varphi_3 - a_{n-1} \varphi'_1 \cdot \varphi_1$$

$${}_n\lambda_3 \cdot \varphi'_2 \cdot {}_n\varphi_3 = 0 = \varphi'_2 A \cdot {}_{n-1}\varphi_3 - \beta_{n-1} \varphi'_2 \cdot \varphi_2$$

whence
$$a_{n-1} = \varphi'_1 A \cdot {}_{n-1}\varphi_3 / \varphi'_1 \cdot \varphi_1$$

$$\beta_{n-1} = \varphi'_2 A \cdot {}_{n-1}\varphi_3 / \varphi'_2 \cdot \varphi_2$$

This determines λ_3, φ_3 , and a similar process with $\gamma\varphi_3, \delta\varphi_4$ etc. can be used to obtain the other latent roots and characteristic vectors.

To check the values of the latent roots use may be made of the relations:

$$\sum_{r=1}^n \lambda_r = \text{tr}|A|$$

$$\lambda_1 \cdot \lambda_2 \cdot \lambda_3 \cdot \dots \cdot \lambda_n = |A|$$

which follow at once from the definition of the characteristic equation given in 7.2.8.

The method just described has the disadvantage that it becomes progressively more complicated as successive latent roots and vectors are determined. A more satisfactory method is the following. Assume that λ_1, φ_1 are respectively the largest latent root and associated latent vector of A , and that φ_1 is normalized so that $\varphi_1' \cdot \varphi_1 = 1$.

Consider next the matrix A_1 where:

$$A_1 = A - \lambda_1 \varphi_1 \cdot \varphi_1' \quad \dots (7.6.4)$$

and the product $\varphi_1 \varphi_1'$ is defined as:

$$\begin{bmatrix} \varphi_1^2 & \varphi_1 \varphi_2 & \varphi_1 \varphi_3 & \dots & \varphi_1 \varphi_n \\ \varphi_2 \varphi_1 & \varphi_2^2 & \varphi_2 \varphi_3 & \dots & \varphi_2 \varphi_n \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \varphi_n \varphi_1 & \varphi_n \varphi_2 & \varphi_n \varphi_3 & \dots & \varphi_n^2 \end{bmatrix}$$

We have, from 7.6.4

$$\begin{aligned} A_1 \varphi_1 &= A \varphi_1 - \lambda_1 \varphi_1 \cdot \varphi_1' \cdot \varphi_1 \\ &= A \varphi_1 - \lambda_1 \varphi_1 \quad (\text{since } \varphi_1 \text{ is normalized and } \varphi_1' \cdot \varphi_1 = 1) \\ &= 0 \end{aligned}$$

Again, if φ_r is any other latent vector of A ,

$$\begin{aligned} A_1 \varphi_r &= A \varphi_r - \lambda_1 \varphi_1 \cdot \varphi_1' \cdot \varphi_r \\ &= A \varphi_r \quad (\text{since the } \varphi_r \text{ is orthogonal to } \varphi_1) \\ &= \lambda_r \varphi_r \end{aligned}$$

Whence A_1 has the same latent vectors as A but the latent root corresponding to φ_1 is now zero.

The iterative method can now be applied to A_1 and will result in the determination of φ_2 and λ_2 . A new matrix, $A_2 = A_1 - \lambda_2 \varphi_2 \cdot \varphi_2'$, is now formed as before, this is easily shown to have the same latent vectors as A but to have latent roots of value zero associated with φ_1 and φ_2 . The iterative method, applied to A_2 , will thus lead to λ_3 and φ_3 . Clearly the process can be repeated and will, in principle, lead to the determination of all of the λ_i and φ_i .

The φ_i are not orthogonal when A is unsymmetric and the above method of removing roots is not available, J. H. Wilkinson has suggested that the following is the simplest procedure in this case. First we define row vectors

$$a_i = (a_{i1} a_{i2} \dots a_{in})$$

such that the matrix A can be written:

$$A \equiv \begin{bmatrix} a_{11} a_{12} \dots a_{1n} \\ a_{21} a_{22} \dots a_{2n} \\ \vdots \\ a_{n1} a_{n2} \dots a_{nn} \end{bmatrix} \equiv \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} \quad \dots (7.6.5)$$

Next, let λ_1 and φ_1 be respectively the largest latent root and associated latent vector of A , assume that φ_{1i} is the largest component of φ_1 . We now form a new vector x_1 which is simply φ_1 so scaled that its largest component φ_{1i} becomes unity, thus:

$$x_1 \equiv \begin{bmatrix} x_{11} \\ x_{12} \\ \vdots \\ x_{1i-1} \\ 1 \\ x_{1i+1} \\ \vdots \\ x_n \end{bmatrix} \equiv \frac{1}{|\varphi_{1i}|} \begin{bmatrix} \varphi_{11} \\ \varphi_{12} \\ \vdots \\ \varphi_{1i} \\ \vdots \\ \varphi_{1n} \end{bmatrix}$$

Now, since latent vectors are undetermined to the extent of a numerical constant,

$$Ax_1 = \lambda_1 x_1 \quad \dots (7.6.6)$$

and, using the notation of 7.6.5,

$$Ax_1 = \begin{bmatrix} a_1 x_1 \\ a_2 x_1 \\ \vdots \\ a_i x_1 \\ \vdots \\ a_n x_1 \end{bmatrix} = \begin{bmatrix} \lambda_1 x_{11} \\ \lambda_1 x_{12} \\ \vdots \\ \lambda_1 x_{1i} \\ \vdots \\ \lambda_1 x_{1n} \end{bmatrix} \quad (\text{from 7.6.6})$$

whence, since by the definition of x_1 , $x_{1i} = 1$,

$$a_i x_1 = \lambda_1 \quad \dots (7.6.7)$$

and similarly, for any other latent root and vector λ_r, x_r ,

$$a_i x_r = \lambda_r \quad \dots (7.6.8)$$

where it is assumed that x_r is so scaled that $x_{ri} = 1$.

We now form the matrix:

$$A_1 = A - x_1 a'_i$$

and notice that, by virtue of 7.6.7,

$$\begin{aligned} A_1 x_1 &= A x_1 - x_1 a'_i x_1 \\ &= \lambda_1 x_1 - x_1 \lambda_1 = 0 \end{aligned}$$

and, using 7.6.8

$$\begin{aligned} A_1 x_r &= A x_r - x_1 a'_i x_r \\ &= \lambda_r x_r - x_1 \lambda_r \\ &= \lambda_r (x_r - x_1) \end{aligned}$$

or, since $A_1 x_1 = 0$, $A_1 (x_r - x_1) = \lambda_r (x_r - x_1)$

We have therefore shown that the matrix A_1 has latent roots $\lambda_2 \dots \lambda_n$ which are identical with those of A but are associated with latent vectors $(x_2 - x_1) (x_3 - x_1) \dots (x_n - x_1)$ where x_1 is the latent root associated with the largest latent vector, λ_1 , of A . The remaining latent root of A_1 is zero.

The iterative method, applied to A_1 , will thus lead to the determination of λ_2 (the sub-dominant latent root of A) and to an associated latent vector $(x_2 - x_1)$.

Now since $(x_2 - x_1)$ is, as usual, undetermined to the extent of a numerical factor, we cannot simply add x_1 (which is already known) to it in order to obtain x_2 . Assume, therefore, that the vector $(x_2 - x_1)$ differs in scale from the x_1 already determined by a factor k . We then have:

$$x_2 = x_1 + k(x_2 - x_1) \quad \dots (7.6.9)$$

where the bracket is to be regarded as an operational symbol and must not be removed. Multiplying 7.6.9 by a_i and using 7.6.7 and 7.6.8 we have

$$\lambda_2 = \lambda_1 + k a_i (x_2 - x_1)$$

whence
$$x_2 = x_1 + \frac{(\lambda_2 - \lambda_1)}{a_i \cdot (x_2 - x_1)} (x_2 - x_1) \quad \dots (7.6.9)$$

It should be noticed that the definition of A_1 implies that its i th row is zero. This, in turn, involves that all of the latent vectors of A_1 have their i th components zero. Thus it is sufficient to work with a matrix derived by striking the i th row and column from A_1 . The root removal process can be continued until all of the λ_r are obtained, a

sequence of matrices A_1, A_2 etc. being derived. The associated latent vectors can be obtained by the application of equations of the type 7.6.9 in sequence starting from the sub-dominant root and vector.

To conclude this section we may mention three other methods which have been used for the calculation of latent roots and characteristic vectors. The first is the so-called 'purification' process of RICHARDSON⁽⁷⁾, which seeks to remove the contribution of φ_i from an arbitrary vector by multiplication by $(A - \lambda_i I)$. This method demands a certain amount of art on the part of the user, and appears unsuitable for use on large systems and with automatic computing machinery. The second method is the so-called 'Escalator' of MORRIS^{(8), (9)} which depends upon the solutions of a characteristic equation of order $(n - 1)$ to obtain the characteristic equation of order (n) .

The operation of the escalator is as follows. Consider the quadratic form associated with a matrix $A_{(n)}$ of order n ; this is:

$$Q_n = \sum_{i,j=1}^n a_{ij} x_i x_j$$

Suppose that we consider the matrix obtained by removing the column a_{in} ($i = 1 \dots n$) and the row a_{nj} ($j = 1 \dots n$) from A ,

$$A_{(n)} = \begin{bmatrix} a_{11} & a_{12} & \dots & a_{1, n-1} & a_{1n} \\ a_{21} & a_{22} & \dots & a_{2, n-1} & a_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots \\ a_{n-1, 1} & a_{n-1, 2} & \dots & a_{n-1, n-1} & a_{n-1, n} \\ a_{n1} & a_{n2} & \dots & a_{nn} & a_{nn} \end{bmatrix}$$

$$= \begin{bmatrix} & & & & a_{1n} \\ & & & & a_{2n} \\ & & & & \vdots \\ & & & & a_{n-1, n} \\ & & & & \\ \dots & \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nn} & a_{nn} \end{bmatrix}$$

Calling this matrix $A_{(n-1)}$ we have the associated quadratic form:

$$Q_{n-1} = \sum_{i,j=1}^{n-1} a_{ij} x_i x_j$$

and we may write:

$$Q_n = Q_{n-1} + x_n \sum_{i=1}^{n-1} a_{in} x_i + x_n \sum_{j=1}^{n-1} a_{nj} x_j + a_{nn} x_n^2 \dots (7.6.10)$$

If we assume that the latent roots λ_r ($r = 1 \dots n-1$) and the characteristic vectors $\varphi_r = (\phi_{r1}, \phi_{r2} \dots \phi_{rn})$ ($r = 1 \dots n-1$) are known for $A_{(n-1)}$, then it is well known, from matrix theory, that we can express Q_{n-1} as:

$$Q_{n-1} = \sum_{r=1}^{n-1} \lambda_r \varphi_r^2 \dots (7.6.11)$$

Now the components ϕ_{rj} of φ_r are in the directions of the original 'axes' x_j so that we have:

$$\varphi_r = l_{r1} x_1 + l_{r2} x_2 \dots + l_{r, n-1} x_{n-1} \quad (r = 1 \dots n-1),$$

and the matrix $L_{(n-1)}$, given by:

$$L_{(n-1)} = \begin{bmatrix} l_{11} & l_{12} & \dots & l_{1, n-1} \\ l_{21} & l_{22} & \dots & l_{2, n-1} \\ \dots & \dots & \dots & \dots \\ l_{n-1, 1} & l_{n-1, 2} & \dots & l_{n-1, n-1} \end{bmatrix}$$

is an ortho-normal one (i.e. it corresponds to a change from one set of orthogonal axes to another set related to the first by a simple rotation and without change of scale). $L_{(n-1)}$ is symmetric and obeys:

$$L_{(n-1)} \cdot L'_{(n-1)} = I.$$

We can thus express the quantities x_r in terms of the φ_r :

$$x_r = l_{1r} \varphi_1 + l_{2r} \varphi_2 \dots + l_{n-1, r} \varphi_{n-1} \quad (r = 1 \dots n-1)$$

so that, using equation 7.6.11, 7.6.10 becomes:

$$Q_n = \sum_{r=1}^{n-1} \lambda_r \varphi_r^2 + x_n \sum_{i=1}^{n-1} a_{in} \sum_{s=1}^{n-1} l_{si} \varphi_s + x_n \sum_{j=1}^{n-1} a_{nj} \sum_{s=1}^{n-1} l_{sj} \varphi_s + a_{nn} x_n^2 \dots (7.6.12)$$

We notice, at this point, that x_n is unchanged by the transformation. It is now evident that the matrix of equation 7.6.12, in terms of $(\varphi_1, \varphi_2 \dots \varphi_{n-1}, x_n)$ is:

$$A_{(n)} = \begin{bmatrix} \lambda_1 & 0 & 0 & \dots & 0 & \sum_{r=1}^{n-1} a_{rn} l_{1r} \\ 0 & \lambda_2 & 0 & \dots & 0 & \sum_{r=1}^{n-1} a_{rn} l_{2r} \\ 0 & 0 & \lambda_3 & \dots & 0 & \dots \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & \dots & \lambda_{n-1} & \sum_{r=1}^{n-1} a_{rn} l_{n-1, r} \\ \sum_{r=1}^{n-1} a_{nr} l_{1r} & \sum_{r=1}^{n-1} a_{nr} l_{2r} & \dots & \sum_{r=1}^{n-1} a_{nr} l_{n-1, r} & a_{nn} \end{bmatrix} \dots (7.6.13)$$

so that the characteristic equation becomes:

$$|A_{(n)} - \lambda I| = \begin{vmatrix} \lambda_1 - \lambda & 0 & 0 & \dots & 0 & a_{rn} l_{1r} \\ 0 & \lambda_2 - \lambda & 0 & \dots & 0 & a_{rn} l_{2r} \\ 0 & 0 & \lambda_3 - \lambda & \dots & 0 & a_{rn} l_{3r} \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \lambda_{n-1} - \lambda & a_{rn} l_{n-1, r} \\ a_{nr} l_{1r} & a_{nr} l_{2r} & a_{nr} l_{3r} & \dots & a_{nr} l_{n-1, r} & a_{nn} - \lambda \end{vmatrix} = 0 \dots (7.6.14)$$

where double suffix notation has been used to avoid writing out the summations.

We now subtract $a_{rn} \cdot l_{1r} / \lambda_1 - \lambda$ times the first column from the last, thus eliminating the first coefficient in column n , then $a_{rn} \cdot l_{2r} / \lambda_2 - \lambda$, times the second column and so on. The result is:

$$\begin{vmatrix} \lambda_1 - \lambda & 0 & 0 & \dots & 0 & 0 \\ 0 & \lambda_2 - \lambda & 0 & \dots & 0 & 0 \\ 0 & 0 & \lambda_3 - \lambda & \dots & 0 & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & \dots & \lambda_{n-1} - \lambda & 0 \\ a_{nr} l_{1r} & a_{nr} l_{2r} & a_{nr} l_{3r} & \dots & a_{nr} l_{n-1, r} & a_{nn} - \lambda - \sum_{i=1}^{n-1} \frac{(a_{nr} l_{ir}) \cdot (a_{rn} l_{ir})}{\lambda_i - \lambda} \end{vmatrix} = 0,$$

whence, expanding:

$$a_{nn} - \lambda = \sum_{i=1}^{n-1} \frac{(a_{nr}l_{ir})(a_{rn}l_{ir})}{\lambda_i - \lambda}$$

or, since A and L are symmetric:

$$a_{nn} - \lambda = \sum_{i=1}^{n-1} \frac{P_i^2}{(\lambda_i - \lambda)} \quad \dots (7.6.15)$$

where

$$P_i = \sum_{r=1}^{n-1} a_{rn}l_{ir} \quad \dots (7.6.16)$$

Notice that P_i is simply the sum of the products of the coefficients of the n th column in the matrix A with the components of the i th characteristic vector in its normalized condition.

Morris has shown that the characteristic vector associated with any solution, λ , of equation 7.6.15 is given by:

$$\frac{\phi_{\lambda,r}}{\phi_{\lambda,n}} = - \sum_{i=1}^{n-1} \frac{l_{ir}P_i}{(\lambda_i - \lambda)} \quad (r = 1 \dots n-1) \quad \dots (7.6.17)$$

where $\phi_{\lambda,r}$ is the r th component of the desired vector and l_{ir} is the r th component of the characteristic vector associated with λ_i , the vector being in its normalized condition. Components, $\phi_{\lambda,r}$, obtained from equation 7.6.17 have to be normalized to give the desired $l_{\lambda,r}$.

The operation of the 'escalator' is as follows:

- (1) Solve $a_{11} - \lambda = 0$ to give λ_1 .
- (2) Find ϕ_1 associated with λ_1 .
- (3) Use equation 7.6.15 to find λ_1 and λ_2 for next higher step.
- (4) Find ϕ_1 and ϕ_2 using λ_1, λ_2 and equation 7.6.17.
- (5) Use λ_1 and λ_2 with ϕ_1 and ϕ_2 to find $\lambda_1, \lambda_2, \lambda_3$ and so on.

The method when applied to an n th order matrix, leads to the intermediate calculation of about $n^2/2$ roots and $n^3/3$ components of characteristic vectors, so that the labour involved for large n is considerable. The method does not appear to be well suited to operation on an automatic calculating machine, but for hand computation for reasonable values of n (up to about 10) it appears very convenient, especially when the checks suggested by Morris are applied.

7.7 DESCENT METHODS FOR THE SOLUTION OF LINEAR EQUATIONS

An alternative class of methods for the solution of sets of linear simultaneous equations depends upon the geometry, in hyperspace, of an associated positive definite quadratic form.

Consider first the expression S (equation 7.5.4). We have:

$$S = \frac{1}{2} \xi \cdot A \xi - b \cdot \xi \quad \dots (7.7.1)$$

and we have seen, in section 7.5, that the surfaces $S = \text{const.}$ represent a set of hyper-ellipsoids whose common centre has co-ordinates $(\xi_1 = x_1, \xi_2 = x_2 \dots \xi_n = x_n)$ which define the solution of the set of simultaneous equations 7.2.1.

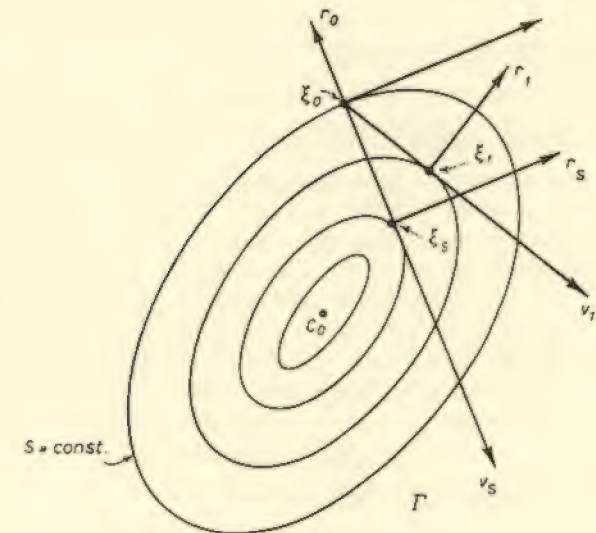


Figure 7.7.1

The problem of obtaining a solution of the original equations may now be reduced to the following equivalent: given any point (ξ_0) , find a process which will lead to a point (ξ_1) for which the value of S is less than that for (ξ_0) . A repetition of such a process will eventually lead to some point (ξ_n) from which no possible decrease in the value of S is possible within the limits of accuracy assigned. The point (ξ_n) is then the 'solution' of the equations.

We shall now show that an infinite number of such processes is possible, and shall discuss some of the more significant variants. First take *any* plane Γ passing through the point (ξ_0) and containing the residue vector (r_0) (see section 7.5). The surfaces $S = \text{const.}$ will intersect this plane in a set of ellipses having a common centre C_0 (Figure 7.7.1).

Consider *any* vector (v_1) lying in Γ and making an angle $> \pi/2$ with the direction of (r_0) , then, for points (ξ_{v_1}) on this vector and sufficiently near to (ξ_0) S will have a smaller value than that at (ξ_0) .

The mode of procedure is clearly to fix upon some (v_1) and proceed along it until S no longer decreases, at (ξ_1) say. To determine (ξ_1) we notice that at (ξ_1) , (v_1) and the residue vector (r_1) are at right angles (notice that r_1 will not generally lie in Γ), and this condition may be ensured by making the vectors v_1 and r_1 have a zero scalar product.

Now any point (ξ) on (v_1) may be written:

$$\xi = \xi_0 + a v_1 \quad \dots (7.7.2)$$

where a is a scalar. Again, by definition:

$$r = A\xi - b \quad (7.7.3)$$

whence

$$r_0 = A\xi_0 - b$$

and consequently:

$$r = r_0 + a A \cdot v_1 \quad \dots (7.7.4)$$

Now at (r_1) we have $(r_1 \cdot v_1) = 0$ whence, from equation 7.7.4 :

$$r_0 \cdot v_1 + a(v_1 \cdot A v_1) = 0$$

or

$$a = - (r_0 \cdot v_1) / (v_1 \cdot A v_1) \quad \dots (7.7.5)$$

so that our end point (ξ_1) is given by:

$$\xi_1 = \xi_0 - \frac{(r_0 \cdot v_1)}{(v_1 \cdot A v_1)} v_1 \quad \dots (7.7.6)$$

So far, the vector (v_1) has been completely arbitrary [except for the general condition of making an angle $> \pi/2$ with the direction of (r_0)], which demonstrates the truth of our original assertion regarding the infinite number of possible descent processes; we now examine some of the more practical aspects of the choice of (v_1) .

Three choices of (v_1) have significant value:

(1) (v_1) is taken parallel to an axis having unit vector (x_k) .

(2) (v_1) is taken in the direction of steepest descent from (ξ_0) —that is at 180° to (r_0) .

(3) (v_1) is so chosen as to pass through C_0 .

The first choice gives the so-called 'relaxation' method for the solution of equations, and if (x_k) is the axis chosen we have:

$$\begin{aligned} \xi_1 &= \xi_0 - \frac{(r_0 \cdot x_k)}{(x_k \cdot A x_k)} \cdot x_k \\ &= \xi_0 - \frac{r_{0k}}{a_{kk}} \cdot x_k. \quad \dots (7.7.7) \end{aligned}$$

That is, the component of ξ_0 in the direction of the x_k th axis is decreased by r_{0k}/a_{kk} . Notice that by r_{0k} we mean the residual from the k th equation at the first approximation.

The virtue of the relaxation method is the simplicity of the formula given in equation 7.7.7. Its defect lies in the fact that x_k has to be chosen by means of some ill-defined criterion, usually by considering the equation which gives the largest residual. In actual computation a skilled operator will often operate on several values of x_k at once, a process known as 'block' or 'group' relaxation. Unfortunately, a high degree of skill is required to make full use of the method in its widest sense, but, even so, the most inexperienced operator can always obtain a result eventually.

For use with an automatic computing machine it is desirable (even if not mandatory) to use a process which calls for no decision based upon the judgement of the user. Such a method is the 'steepest descent,' for which we take $v_1 = -r_0$, and thus obtain:

$$\begin{aligned} \xi_1 &= \xi_0 - \frac{(r_0 \cdot r_0)}{(r_0 \cdot A r_0)} \cdot r_0 \quad \dots (7.7.8) \\ &= \xi_0 - \left(\frac{\sum_{i=1}^n r_{0i}^2}{\sum_{i,j=1}^n a_{ij} r_{0i} \cdot r_{0j}} \right) \cdot r_0 \end{aligned}$$

(r_{0i} is the i th component of r_0).

For well-conditioned equations this method is excellent; when, however, it is applied to an ill-conditioned set, the situation shown in Figure 7.7.2 may arise, in which the minimum ξ_1 slightly overshoots a principal diameter of the ellipse system.

Under these circumstances the solution may oscillate from side to side of the principal diameter and, in consequence, require many

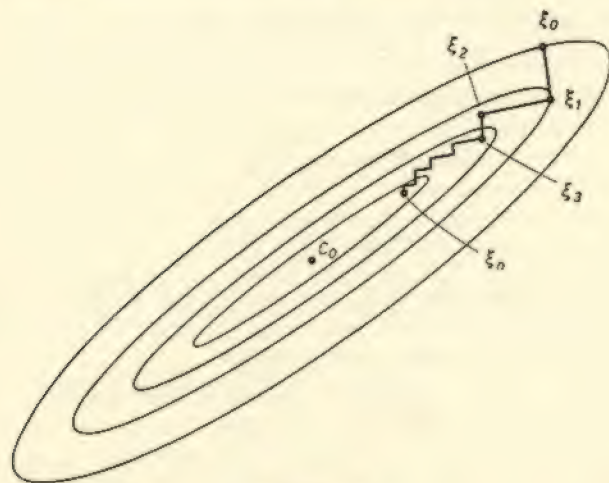


Figure 7.7.2

iterations to reach C_0 . Experience has shown that more rapid convergence may be obtained, for ill-conditioned systems, by taking:

$$\xi_1 = \xi_0 - \frac{(\mathbf{r}_0 \cdot \mathbf{r}_0)}{(\mathbf{r}_0 \cdot A\mathbf{r}_0)} \cdot \mathbf{r}_0 \quad \dots (7.7.9)$$

for, say, four iterations and then making a single application of the exact formula 7.7.8.

The third method of procedure may be derived as follows. Suppose that, at the previous stage of the process, there was derived a vector \mathbf{v}_n and a point on it, ξ_n , such that $(\mathbf{r}_n \cdot \mathbf{v}_n) = 0$, just as was done in the previous derivations. We now have a pair of orthogonal vectors \mathbf{v}_n and \mathbf{r}_n . Consider the system of ellipses $S = \text{const.}$ in the plane of \mathbf{v}_n and \mathbf{r}_n shown in Figure 7.7.3.

Now any vector in the plane Γ containing \mathbf{r}_n and \mathbf{v}_n can be written:

$$\mathbf{v}_{n+1} = a\mathbf{r}_n + \beta\mathbf{v}_n. \quad \dots (7.7.10)$$

If C_{n+1} is the centre of the ellipses $S = \text{const.}$ then, at C_{n+1} , the residue vector \mathbf{r}_{n+1} is perpendicular to Γ , and consequently to both \mathbf{r}_n and \mathbf{v}_n , so that $(\mathbf{r}_{n+1} \cdot \mathbf{r}_n) = (\mathbf{r}_{n+1} \cdot \mathbf{v}_n) = 0$.

Again:

$$\xi_{n+1} = \xi_n + \mathbf{v}_{n+1} = \xi_n + a\mathbf{r}_n + \beta\mathbf{v}_n. \quad \dots (7.7.11)$$

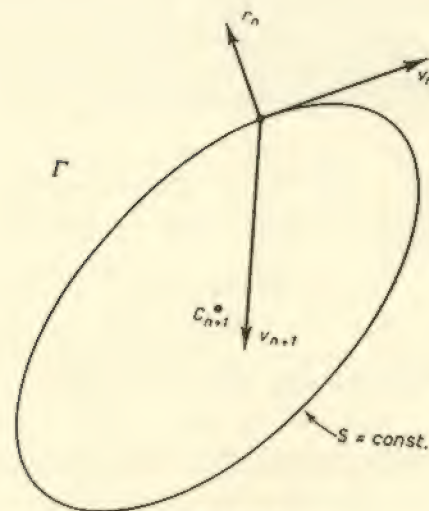


Figure 7.7.3

And since

$$\mathbf{r} = A\xi - \mathbf{b}.$$

$$\mathbf{r}_{n+1} = \mathbf{r}_n + A(\xi_{n+1} - \xi_n) = \mathbf{r}_n + aA\mathbf{r}_n + \beta A\mathbf{v}_n. \quad \dots (7.7.12)$$

But, from the orthogonality conditions:

$$(\mathbf{r}_n \cdot \mathbf{r}_{n+1}) = 0 = (\mathbf{r}_n \cdot \mathbf{r}_n) + a(\mathbf{r}_n \cdot A\mathbf{r}_n) + \beta(\mathbf{r}_n \cdot A\mathbf{v}_n)$$

$$(\mathbf{v}_n \cdot \mathbf{r}_{n+1}) = 0 = a(\mathbf{v}_n \cdot A\mathbf{r}_n) + \beta(\mathbf{v}_n \cdot A\mathbf{v}_n)$$

(\mathbf{v}_n and \mathbf{r}_n are *ab initio* orthogonal so that $(\mathbf{r}_n \cdot \mathbf{v}_n) = 0$)

whence:

$$\left. \begin{aligned} a &= \frac{(\mathbf{r}_n \cdot \mathbf{r}_n)(\mathbf{v}_n \cdot A\mathbf{v}_n)}{(\mathbf{v}_n \cdot A\mathbf{r}_n)^2 - (\mathbf{r}_n \cdot A\mathbf{r}_n)(\mathbf{v}_n \cdot A\mathbf{v}_n)} \\ \beta &= \frac{-(\mathbf{r}_n \cdot \mathbf{r}_n)(\mathbf{v}_n \cdot A\mathbf{r}_n)}{(\mathbf{v}_n \cdot A\mathbf{r}_n)^2 - (\mathbf{r}_n \cdot A\mathbf{r}_n)(\mathbf{v}_n \cdot A\mathbf{v}_n)} \end{aligned} \right\} \quad \dots (7.7.13)$$

The set of formulae (7.7.11–7.7.13) may be regarded as the two-dimensional analogue of the steepest descent formula (7.7.8); they have the advantage that the solution is obtained in exactly the same number of steps as the equations have variables. This is due to the fact that each new steepest descent vector is orthogonal to all previous descent vectors. Thus, since the hyperspace has only n dimensions, only n such vectors exist and the process is completed at the n th vector. In this respect the method is closely related to the Gram–Schmidt orthogonalization procedure.

7.8 A CAUTIONARY PARAGRAPH

We have mentioned, in various places, the fact that the residue criterion based upon:

$$S = \frac{1}{2} \xi \cdot A \xi - \mathbf{b} \cdot \xi$$

assumes that the matrix A is both symmetric and positive-definite. It is worth reminding the potential user of 7.7.6–7.7.13 that these equations must only be used when A is known to obey these conditions.

A cautionary example is provided by the equations:

$$\begin{aligned} x_1 - 2x_2 + 3x_3 + x_4 &= 3 \\ -2x_1 + x_2 - 2x_3 - x_4 &= -4 \\ 3x_1 - 2x_2 + x_3 + 5x_4 &= 7 \\ x_1 - x_2 + 5x_3 + 3x_4 &= 8 \end{aligned}$$

for which the solution is (1, 1, 1, 1). An application of equation 7.7.8 from the starting point ξ_0 (0,0,0,0) for which

$$\mathbf{r}_0 = (-3, +4, -7, -8)$$

leads to

$$\xi_1 = (+.34, -.45, +.79, +.90)$$

with

$$\mathbf{r}_1 = (+1.51, +.39, +.21, -.56)$$

which is quite satisfactory [$(\mathbf{r}_0 \cdot \mathbf{r}_0) = 138$, $(\mathbf{r}_1 \cdot \mathbf{r}_1) = 2.79$]. A further application of the procedure leads, however, to co-ordinate changes of large magnitude $O(10)$ and to an increase in $(\mathbf{r}_2 \cdot \mathbf{r}_2)$ to 1475.49!

The reason for this behaviour is clear, since when A is not positive-definite, the expression S represents, in general, hyper-elliptic-hyperboloids, and the geometrical assumptions made in section 7.7 no longer hold.

7.9 DESCENT METHODS BASED ON R^2

By means of the positive definite form:

$$R^2 = (\mathbf{r} \cdot \mathbf{r}) = (A\xi - \mathbf{b})(A\xi - \mathbf{b}) \quad \dots (7.9.1)$$

we may construct the analogues of the descent formulae previously derived for S , but with the advantage that, since R^2 is known to be positive-definite, this limitation on A is removed. As a penalty for this the formulae are, of course, slightly more complicated.

We again choose an arbitrary vector (\mathbf{v}_1) and seek to find the least value of R^2 when some multiple ($a\mathbf{v}_1$) is added to ξ . Thus:

$$\begin{aligned} R_{(\xi+a\mathbf{v}_1)}^2 &= [A(\xi + a\mathbf{v}_1) - \mathbf{b}] \cdot [A(\xi + a\mathbf{v}_1) - \mathbf{b}] \\ &= (A\xi \cdot A\xi) + a^2(A\mathbf{v}_1 \cdot A\mathbf{v}_1) + 2a[(A\xi \cdot A\mathbf{v}_1) - (\mathbf{b} \cdot A\mathbf{v}_1)] \\ &\quad - 2(\mathbf{b} \cdot A\xi) + (\mathbf{b} \cdot \mathbf{b}). \end{aligned}$$

At the minimum we have

$$\frac{dR^2}{da} = 0$$

$$\text{whence: } a(A\mathbf{v}_1 \cdot A\mathbf{v}_1) + [A\mathbf{v}_1 \cdot (A\xi - \mathbf{b})] = 0$$

$$\text{or } a = - (A\mathbf{v}_1 \cdot \mathbf{r}_0) / (A\mathbf{v}_1 \cdot A\mathbf{v}_1)$$

which is easily shown to give a *minimum* value for R^2 . It follows that the correction formula corresponding to equation 7.7.6 is:

$$\xi_1 = \xi_0 - \frac{(\mathbf{r}_0 \cdot A\mathbf{v}_1)}{(A\mathbf{v}_1 \cdot A\mathbf{v}_1)} \cdot \mathbf{v}_1 \quad \dots (7.9.2)$$

Just as in the previous case we may choose \mathbf{v}_1 in a number of ways, if it is taken parallel to an axis having unit vector \mathbf{x}_k , we obtain the relaxation formula:

$$\xi_1 = \xi_0 - \frac{(\mathbf{r}_0 \cdot A\mathbf{x}_k)}{(A\mathbf{x}_k \cdot A\mathbf{x}_k)} \cdot \mathbf{x}_k \quad \dots (7.9.3)$$

or, in terms of the matrix elements

$$\xi_1 = \xi_0 - \frac{\sum_{i=1}^n r_{0i} a_{ik}}{\sum_{i=1}^n a_{ik}^2} \cdot \mathbf{x}_k.$$

The steepest descent method is obtained by choosing (\mathbf{v}_1) to have the direction in which R^2 , in the region of ξ_0 , is changing most rapidly. This is, in vector notation:

$$\mathbf{v}_1 = - \text{grad } R^2 \quad \dots (7.9.4)$$

where the k th component of \mathbf{v}_1 , v_{1k} , is given by:

$$v_{1k} = -\frac{\partial}{\partial \xi_k} (R^2). \quad \dots (7.9.5)$$

Now an inspection of the quadratic form of R^2 , given in equation 7.5.3, shows that:

$$\begin{aligned} \frac{\partial}{\partial \xi_k} (R^2) &= \sum_{i=1}^n 2a_{ik} \sum_{j=1}^n (a_{ij}\xi_j - b_i) \\ &= \sum_{i=1}^n 2a_{ik}r_i \end{aligned} \quad \dots (7.9.6)$$

whence

$$\text{grad } (R^2) = 2A'r \quad \dots (7.9.7)$$

where A' is the transpose of A and \mathbf{r} is the residue vector. We thus obtain

$$\xi_1 = \xi_0 - \frac{(\mathbf{r}_0 \cdot AA'\mathbf{r}_0)}{[(AA'\mathbf{r}_0) \cdot (AA'\mathbf{r}_0)]} \cdot A'\mathbf{r}_0 \quad \dots (7.9.8)$$

which is the analogue of equation 7.7.8.

Finally, we may find the minimum of R^2 over a two dimensional plane containing any pair of non-parallel vectors. It is convenient to choose these to be $(\text{grad } R^2)$ at ξ_n and the vector \mathbf{v}_n used to derive ξ_n from ξ_{n-1} —this pair being orthogonal. Thus:

$$\xi_{n+1} = \xi_n + \mathbf{v}_{n+1} \quad \dots (7.9.9)$$

$$\mathbf{v}_{n+1} = \alpha(\text{grad } R^2)_n + \beta \mathbf{v}_n = 2\alpha A'\mathbf{r}_n + \beta \mathbf{v}_n \quad \dots (7.9.10)$$

where α and β are determined so as to make $R^2(\xi_{n+1})$ a minimum. Using 7.9.9, 7.9.10 and the definition of R^2 7.9.1, we have:

$$\begin{aligned} R^2(\xi_{n+1}) &= (A\xi_{n+1} - \mathbf{b}) \cdot (A\xi_{n+1} - \mathbf{b}) \\ &= [A\xi_n - \mathbf{b} + \alpha A(\text{grad } R^2)_n + \beta A\mathbf{v}_n] \\ &\quad [A\xi_n - \mathbf{b} + \alpha A(\text{grad } R^2)_n + \beta A\mathbf{v}_n] \\ &= (\mathbf{r}_n + 2\alpha AA'\mathbf{r}_n + \beta A\mathbf{v}_n) \cdot (\mathbf{r}_n + 2\alpha AA'\mathbf{r}_n + \beta A\mathbf{v}_n). \end{aligned}$$

Now we require, for a minimum:

$$\frac{\partial R^2}{\partial \alpha} = \frac{\partial R^2}{\partial \beta} = 0$$

whence:

$$4AA'\mathbf{r}_n \cdot (\mathbf{r}_n + 2\alpha AA'\mathbf{r}_n + \beta A\mathbf{v}_n) = 0$$

$$2A\mathbf{v}_n \cdot (\mathbf{r}_n + 2\alpha AA'\mathbf{r}_n + \beta A\mathbf{v}_n) = 0$$

giving:

$$\left. \begin{aligned} 2\alpha &= \frac{(\mathbf{r}_n \cdot A\mathbf{v}_n)(A\mathbf{v}_n \cdot AA'\mathbf{r}_n) - (\mathbf{r}_n \cdot AA'\mathbf{r}_n)(A\mathbf{v}_n \cdot A\mathbf{v}_n)}{(AA'\mathbf{r}_n \cdot AA'\mathbf{r}_n)(A\mathbf{v}_n \cdot A\mathbf{v}_n) - (A\mathbf{v}_n \cdot AA'\mathbf{r}_n)(A\mathbf{v}_n \cdot AA'\mathbf{r}_n)} \\ \beta &= \frac{(\mathbf{r}_n \cdot AA'\mathbf{r}_n)(A\mathbf{v}_n \cdot AA'\mathbf{r}_n) - (\mathbf{r}_n \cdot A\mathbf{v}_n)(AA'\mathbf{r}_n \cdot AA'\mathbf{r}_n)}{(AA'\mathbf{r}_n \cdot AA'\mathbf{r}_n)(A\mathbf{v}_n \cdot A\mathbf{v}_n) - (A\mathbf{v}_n \cdot AA'\mathbf{r}_n)(A\mathbf{v}_n \cdot AA'\mathbf{r}_n)} \end{aligned} \right\} \dots (7.9.11)$$

These formulae are far too cumbersome for use by a human computer, but may possibly find application with an automatic digital calculator.

7.10 A NUMERICAL EXAMPLE

To illustrate the merits (and otherwise) of the preceding methods, we give their application to the set of equations proposed in section 7.8. It has already been pointed out that methods based upon S will fail for this set of equations since A is not positive-definite.

Using the relaxation technique of section 7.9, based on R^2 and the initial vector $\xi_0 = (0, 0, 0, 0)$, we obtain *Table 7.10.1* by a direct application of equation 7.9.3.

Table 7.10.1 Relaxation on set of equations of section 7.8

k	ξ_0	r_0	ξ_1	r_1	ξ_2	r_2	ξ_3	r_3	ξ_4	r_4
1	0	-3	0	-1.17	0	+ .231	0	- .633	0	- .135
2	0	+4	0	+2.17	0	+1.235	.432	+1.667	.432	+1.335
3	0	-7	0	+2.17	.466	+2.631	.466	+1.767	.632	+1.933
4	0	-8	1.83	-2.50	1.833	- .171	1.833	- .603	1.833	+ .227
	$(\mathbf{r}_0 \cdot \mathbf{r}_0) = 138$		$(\mathbf{r}_1 \cdot \mathbf{r}_1) = 17.00$		$(\mathbf{r}_2 \cdot \mathbf{r}_2) = 8.53$		$(\mathbf{r}_3 \cdot \mathbf{r}_3) = 6.67$		$(\mathbf{r}_4 \cdot \mathbf{r}_4) = 5.59$	

The residues which indicate the value of k (i.e. the co-ordinate number) in equation 7.9.3 have been underlined. We notice that it is not always possible to select the co-ordinate of largest residual at each step since, as happens in $k = 4$ of \mathbf{r} , a relaxation with respect

SIMULTANEOUS LINEAR EQUATIONS

to this vector has just been performed so that no further improvement can be produced by further alteration to x_4 at this stage.

Next we give the results of the same number of steepest descents applied to the same set of equations.

Table 7.10.2. Steepest descent

$\xi, r_0 (A'r_0)$	ξ_1 r_1 $(A'r_1)$	ξ_2 r_2 $(A'r_2)$	ξ_3 r_3 $(A'r_3)$	ξ_4 r_4
0 -3 -40	+ .488 +1.416 +.230	+ .392 +.488 +.605	+ .384 +.366 -.603	+ .426 +.478
0 +4 +32	-.390 +.267 -1.863	+ .391 +.612 -1.292	+ .409 +.709 -.290	+ .429 +.578
0 -7 -64	+ .781 +.050 -.246	+ .884 +.413 +1.163	+ .868 +.182 -.623	+ .911 +.216
0 -8 -66	+ .805 -.802 -1.007	+1.227 +.102 +2.247	+1.196 -.097 +.276	+1.177 +.083
$(r_0, r_0) = 138$	$(r_1, r_1) = 2.72$	$(r_2, r_2) = .79$	$(r_3, r_3) = .68$	$(r_4, r_4) = .62$
$\xi_1 =$	$\xi_2 =$	$\xi_3 =$	$\xi_4 =$	
$\xi_0 - .0122(A'r_0)$	$\xi_1 + .419(A'r_1)$	$\xi_2 + .0136(A'r_2)$	$\xi_3 + .069(A'r_3)$	

It will be noticed that, initially, the convergence of the steepest descent process is faster than that of the relaxation, but that from ξ_3 there is an oscillation of the type indicated in Figure 7.7.2. This is to be expected in the present set of equations which are quite ill-conditioned.

Finally, Table 7.10.3 shows the results of four descent to ellipse centre applications based on equations 7.9.9-7.9.11. Since ξ_1 is the same as a normal steepest descent and is given in Table 7.10.2 and since $v_1 = (A'r_0)$ and r_1 are also given in the latter table, we do not reproduce them here.

Table 7.10.3

$$(A'r_1) = \begin{bmatrix} +.230 \\ -1.863 \\ -.246 \\ -1.007 \end{bmatrix} \quad (AA'r_1) = \begin{bmatrix} +2.211 \\ -.824 \\ -.865 \\ -2.158 \end{bmatrix} \quad Av_1 (= AA'r_0) = \begin{bmatrix} -362 \\ +306 \\ -578 \\ -590 \end{bmatrix}$$

$$(r_1, AA'r_1) = +4.60, \quad (AA'r_1, AA'r_1) = +10.973,$$

$$(r_1, Av_1) = +13.39, \quad (Av_1, AA'r_1) = +720.67,$$

$$(Av_1, Av_1) = +906864$$

whence:

$$2a = -.441, \beta = +.000336$$

$$v_2 = \begin{bmatrix} -.114 \\ +.833 \\ +.086 \\ +.422 \end{bmatrix} \quad \xi_2 = \begin{bmatrix} +.374 \\ +.443 \\ +.867 \\ +1.227 \end{bmatrix} \quad r_2 = \begin{bmatrix} +.316 \\ +.734 \\ +.238 \\ -.053 \end{bmatrix} \quad (r_2, r_2) = .70.$$

MATRIX 'ROTATION' AND LATENT ROOTS

In a similar manner we obtain:

$$v_3 = \begin{bmatrix} +.080 \\ +.136 \\ +.109 \\ -.074 \end{bmatrix} \quad \xi_3 = \begin{bmatrix} +.454 \\ +.579 \\ +.976 \\ +1.153 \end{bmatrix} \quad r_3 = \begin{bmatrix} +.377 \\ +.566 \\ -.055 \\ +.214 \end{bmatrix} \quad (r_3, r_3) = .51$$

and:

$$v_4 = \begin{bmatrix} +.097 \\ +.102 \\ -.005 \\ -.057 \end{bmatrix} \quad \xi_4 = \begin{bmatrix} +.551 \\ +.681 \\ +.971 \\ +1.096 \end{bmatrix} \quad r_4 = \begin{bmatrix} +.198 \\ +.541 \\ -.258 \\ +.013 \end{bmatrix} \quad (r_4, r_4) = .40.$$

7.11 MATRIX 'ROTATION' AND LATENT ROOTS

Probably the simplest methods of determining both latent roots and latent vectors of a matrix are those in which the original matrix is transformed into some form from which the roots and vectors can be derived directly. The first such method is due to Jacobi and was re-discovered by VON NEUMANN and GOLDSTINE⁽¹⁰⁾ in the late 1940s. In this method the matrix is transformed into truly diagonal form by means of an orthogonal transformation matrix.

Suppose that A is the matrix whose latent roots and vectors are required and that T is an orthogonal transformation matrix so that:

$$TT' = T'T = I \quad \dots (7.11.1)$$

Now, $T'(A - \lambda I) = (T'AT - \lambda I)$ from 7.11.1 whence the λ_i , which are the roots of $|A - \lambda| = 0$, are also the roots of $|T'AT - \lambda| = 0$, that is, the latent roots of $T'AT$ are also those of A .

Now let (λ_i, y_i) be respectively a latent root and associated latent vector of $T'AT$. Then by definition,

$$T'ATy_i = \lambda_i y_i$$

or, pre-multiplying by T , and using 7.11.1,

$$ATy_i = \lambda_i Ty_i$$

whence

$$(A - \lambda_i I)Ty_i = 0$$

which shows that Ty_i is the latent root of A which corresponds to λ_i . The Jacobi method produces a matrix T which has the property that:

$$T'AT \equiv D \equiv \begin{bmatrix} d_{11} & & & 0 \\ & d_{22} & & \\ & 0 & \ddots & \\ & & & d_{nn} \end{bmatrix}$$

from which, by inspection,

$$\lambda_1 = d_{11}, y_1 = (1, 0, 0, 0, \dots, 0); \lambda_2 = d_{22}, y_2 = (0, 1, 0, 0, \dots, 0); \\ \dots \lambda_i = d_{ii} \text{ etc.}$$

T is determined as the product of a sequence of elementary transformations, C_n , each of which produces a matrix in which the largest, off-diagonal, element of the preceding matrix is reduced to zero.

The basis of the method is the fact that (see section 7.5) the latent roots of the matrix A are the squares of the inverse principal axes of the quadratic form:

$$\sum_{i,j=1}^n a_{ij} x_i x_j = 1 \quad (a_{ij} = a_{ji}) \quad \dots (7.11.2)$$

The orthonormal matrix T has the property that:

$$z = T.x$$

transforms 7.11.2 into

$$\sum_{i=1}^n \lambda_i z_i^2 = 1 \quad \dots (7.11.3)$$

where the λ_i are the latent roots of A . It is well known⁽¹¹⁾ that if C is any orthonormal transformation matrix so that the transformed variables, y are given by

$$y = C.x \quad (C'.C = I)$$

then the matrix of the transformed quadratic form in y , corresponding to 7.11.2, is $C'.A.C$. It should be noted that the sum of the squares of the diagonal elements of C is greatest when $C = T$, that is, when the transformed matrix of coefficients is diagonal.

The Jacobi process determines T as the limit of a sequence of simple rotations involving only two of the axes (x_i) at a time, the idea being to eliminate the largest (a_{ij}) from A ($i \neq j$) and then the next largest, and so on.

Suppose that a_{rs} ($= a_{sr}$) is the largest off-diagonal element of A , and that we seek to eliminate it by a rotation of axes through an

angle θ_1 in the plane of (x_r, x_s). The equations of transformation are:

$$\left. \begin{aligned} x_r &= \cos \theta_1 \cdot y_r - \sin \theta_1 \cdot y_s \\ x_s &= \sin \theta_1 \cdot y_r + \cos \theta_1 \cdot y_s \\ x_i &= y_i \quad (i \neq r, s) \end{aligned} \right\} \quad \dots (7.11.4)$$

and substitution in equation 7.11.2 gives for the coefficient of (y_r, y_s)

$$2 \cos \theta_1 \sin \theta_1 (a_{rr} - a_{ss}) - 2 (\cos^2 \theta_1 - \sin^2 \theta_1) a_{rs}$$

so that to eliminate this we take:

$$\tan 2\theta_1 = 2a_{rs}/(a_{rr} - a_{ss}) \quad \dots (7.11.5)$$

which defines the angle of rotation (θ_1).

Now the transformation matrix C_1 say, corresponding to equation 7.11.4 is:

$$C_1 = \begin{bmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 0 & 1 & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \cos \theta_1 & 0 & 0 & \dots & 0 & -\sin \theta_1 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & 1 & 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & \dots & \sin \theta_1 & 0 & 0 & \dots & 0 & \cos \theta_1 & \dots & 0 \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots & \ddots \\ 0 & 0 & 0 & 0 & \dots & 0 & 0 & 0 & \dots & 0 & 0 & 0 & 1 \end{bmatrix} \quad \dots (7.11.6)$$

$\uparrow \qquad \qquad \qquad \uparrow$
 $r\text{th column} \qquad \qquad s\text{th column}$

and it is easily proved, by direct substitution, that the sum of the squares of the elements on the diagonal of the transformed matrix:

$$A_1 = C_1'AC_1$$

is simply:

$$2a_{rs}^2 + \sum_{i=1}^n a_{ii}^2 \quad \dots (7.11.7)$$

which shows, since it is greater than $\sum_{i=1}^n a_{ii}^2$, that A_1 is nearer (in some sense) to the required diagonal matrix than was A .

A repetition of the above process, with the largest off-diagonal element of A_1 will lead to a new transformed matrix A_2 , thus:

$$\begin{aligned} A_1 &= C_1' A C_1 \\ A_2 &= C_2' A_1 C_2 \\ A_3 &= C_3' A_2 C_3 \\ &\dots \dots \dots \\ A_m &= C_m' A_{m-1} C_m \\ &\text{etc.} \end{aligned}$$

where A_m tends to the matrix (λ_i) and the continued product $C_1 C_2 C_3 \dots$ tends to the diagonalizing matrix T defined in equation 7.11.1.

The Jacobi method has the advantage that all of the operations involved are very simple. Convergence is rapid in suitable cases but it must be remembered that, in principle, the method is one of successive approximation and not a finite step one. A minor disadvantage, from the viewpoint of the automatic digital computer, is the fact that the off-diagonal elements must be scanned at each stage to determine which is greatest.

Starting with a suggestion of LANCZOS⁽¹²⁾ in 1950, new methods of latent root and vector determination have been developed in which the matrix is not diagonalized but is reduced to triple-diagonal form. It turns out that transformation to this form can be made in a finite number of steps and that the latent roots and vectors of the triple-diagonal matrix can be fairly easily determined. The two best-known methods of transformation are those of GIVENS⁽¹³⁾ and HOUSEHOLDER⁽¹⁴⁾.

In Givens' method co-ordinate transformations are again applied but this time so that the element reduced to zero is $a_{r-1,s}$ ($r > 1$) and the transformation, corresponding to 7.11.4, has:

$$\tan \theta_1 = a_{r-1,s} / a_{r-1,r}$$

It is no longer necessary to determine the largest off-diagonal element and the elimination proceeds systematically with the liquidation of

$$\begin{array}{ccccccc} a_{13}, & a_{14}, & a_{15}, & \dots & a_{1n} \\ & a_{24}, & a_{25}, & \dots & a_{2n} \\ & & \dots & \dots & \dots \\ & & & & a_{n-2,n} \end{array}$$

The method has the advantage, over that of Jacobi, that the elements reduced to zero remain null so that the labour of transformation decreases as the work proceeds. It has been estimated by

WILKINSON⁽¹⁵⁾ that the work of transformation to triple original form is only about one-twentieth as long as that involved in an equivalent reduction to diagonal form by Jacobi's method.

The second method of reduction is that of Householder. We define $T = T_{n-1} T_{n-3} \dots T_2$ for an n th order matrix A . Here

$$T_m = I - 2w_m w_m'$$

where

$$w_m w_m' = 1$$

and

$$w_m' = (0, 0, \dots, 0, x_m, x_{m+1}, \dots, x_n)$$

whence

$$\sum_{i=m}^n x_i^2 = 1$$

Starting with T_2 the matrix:

$$A^2 = T_2' A T_2$$

is formed with such a choice of x_2, x_3, \dots, x_n that the elements $a_{13} \dots a_{1n}$ of the first row of A are reduced to zero in $A^{(2)}$. Next $A^{(3)}$ is found from $A^{(3)} = T_3' A^{(2)} T_3$ with a choice of the $x_3 \dots x_n$ to make the elements $a_{24} \dots a_{2n}$ zero in $A^{(3)}$. The process continues, with increasing ease, until $A^{(n-1)}$ has been found. This is in triple diagonal form.

The defining equations, the requirement for zero coefficients in the transformed matrix, and the consideration that the sums of squares of elements in any row of the transformed matrix must be equal to that before transformation are adequate to evaluate the x_i at each stage with the exception of a sign which is so chosen as to ensure maximum numerical stability.

WILKINSON⁽¹⁶⁾ has shown that, for the 4×4 matrix:

$$A = \begin{bmatrix} a_1 & b_1 & c_1 & d_1 \\ b_1 & b_2 & c_2 & d_2 \\ c_1 & c_2 & c_3 & d_3 \\ d_1 & d_2 & d_3 & d_4 \end{bmatrix}$$

the optimum choice of values for $(0, x_2, x_3, x_4)$ is

$$x_2^2 = \frac{1}{2} \left[1 + \frac{b_1 s(b_1)}{\sqrt{s}} \right], \quad x_3^2 = \frac{c_1 s(b_1)}{2x_2 \sqrt{s}}, \quad x_4^2 = \frac{d_1 s(b_1)}{2x_2 \sqrt{s}}$$

where $s = b_1^2 + c_1^2 + d_1^2$ and $s(b_1)$ means 'sign of b_1 '.

The number of multiplications required in Householder's method for reduction to triple-diagonal form is about $\frac{2}{3}n^3$ compared with

$\frac{4}{3}n^3$ in Givens' method and $2n^3$ in Lanczos'. Furthermore, Givens' method requires the extraction of $\frac{1}{2}n^2$ square roots compared with the $2n$ of the Householder process. Both methods can be extended to cover the case of unsymmetric matrices although here the reduction leads to almost triangular rather than tri-diagonal form.

The last phase of the determination of latent roots and vectors by the Givens or Householder methods is the determination of the actual λ_i from the triple-diagonal matrix. Several possible methods for doing this exist⁽¹⁷⁾ but the simplest conceptually and the most stable numerically, is the following.

First note that the matrix is assumed to be:

$$T'AT = \begin{bmatrix} \alpha_1 & \beta_2 & 0 & 0 & . & . & . & 0 \\ \beta_2 & \alpha_2 & \beta_3 & 0 & . & . & . & 0 \\ 0 & \beta_3 & \alpha_3 & \beta_4 & . & . & . & 0 \\ . & . & . & . & . & . & . & . \\ 0 & 0 & . & . & . & 0 & \beta_n & \alpha_n \end{bmatrix}$$

The leading principal minors, p_r , are:

$$p_0 = 1, \quad p_1 = \alpha_1, \quad p_2 = \begin{vmatrix} \alpha_1 & \beta_2 \\ \beta_2 & \alpha_2 \end{vmatrix}, \quad p_3 = \begin{vmatrix} \alpha_1 & \beta_2 & 0 \\ \beta_2 & \alpha_2 & \beta_3 \\ 0 & \beta_3 & \alpha_3 \end{vmatrix} \text{ etc.}$$

Denoting by $p_r(\lambda)$ the r th leading principal minor of $(T'AT - \lambda I)$ it is easy to see that:

$$\begin{aligned} p_0(\lambda) &= 1 \\ p_1(\lambda) &= \alpha_1 - \lambda \\ p_2(\lambda) &= (\alpha_1 - \lambda)(\alpha_2 - \lambda) - \beta_2^2 \\ &\dots \dots \dots \\ p_r(\lambda) &= (\alpha_r - \lambda)p_{r-1}(\lambda) - \beta_r^2 p_{r-2}(\lambda) \end{aligned}$$

It can be shown that, if when any $p_r(\lambda)$ is zero its sign is assumed to be the opposite of that of $p_{r-1}(\lambda)$, the number of agreements in sign, σ_λ , between consecutive numbers of the sequence $p_0(\lambda)$, $p_1(\lambda)$, $p_2(\lambda)$, \dots , $p_n(\lambda)$ is equal to the number of latent roots of $T'AT$ which are greater than λ .

To use the method a pair of numbers (a, b) is found such that $\sigma_a \geq K > \sigma_b$. The value of $\sigma_{\frac{a+b}{2}}$ is then computed. If $\sigma_a > \sigma_{\frac{a+b}{2}}$ then a root, λ , lies between a and $\frac{a+b}{2}$; if, however, $\sigma_{\frac{a+b}{2}} > \sigma_b$ the

root lies between $\frac{a+b}{2}$ and b . According to which criterion is satisfied, either $\sigma_{\frac{1}{2}(a+\frac{a+b}{2})}$ or $\sigma_{\frac{1}{2}(\frac{a+b}{2}+b)}$ is computed and the test again applied. It is seen that the precision of location is increased by a factor of 2 at each partition and that the convergence is independent of root separation and magnitude.

To start the process several strategies exist. For a complete determination of roots, a determination of λ_{\max} by the iterative method followed by a survey of the values of the $p_r(\lambda)$ for lesser values of λ will give the rough positions of roots and this can be followed by a precise determination of any desired roots by the method just outlined.

On the other hand, a rough bound for λ_{\max} may be found as follows. If $x = (x_1, x_2, \dots, x_n)$ is any latent vector of $T'AT$ and λ is its associated root, then, since

$$T'ATx = \lambda x$$

$$\beta_r x_{r-1} + \alpha_r x_r + \beta_{r+1} x_{r+1} = \lambda x_r$$

Assume that x_r is the numerically greatest component of x . We then have the inequality:

$$|\lambda| \leq |\beta_r| + |\alpha_r| + |\beta_{r+1}|$$

and a survey of $p_r(\lambda)$ for values of λ between $-(|\beta_r| + |\alpha_r| + |\beta_{r+1}|)$ and $+(|\beta_r| + |\alpha_r| + |\beta_{r+1}|)$ must entrap all of the latent roots.

The accurate determination of the associated latent vectors, once the latent roots are known, is still a non-trivial problem. Wilkinson suggests that the best method is as follows. For the latent root λ_i calculate the matrix $(T'AT - \lambda_i I)$. Reduce this to an upper triangular matrix U_i (say) and then solve the equations $U_i x = e$, where $e' = (1, 1, 1, \dots, 1)$. The solution x_i is then the required latent vector.

7.12 MONTE CARLO METHODS

The somewhat intriguing title of this section is now applied to methods of numerical analysis which make use of the theory of Games⁽¹⁸⁾. The suggestion that such methods might be applicable appears to have been due, originally, to von Neumann and Ulam; but statisticians have since claimed that the techniques are simply those of small sample analysis and, as such, have been known for a considerable time.

Monte Carlo methods appear to be applicable to most problems involving *linear* operators, and the present example of a technique for determining the elements of an inverse matrix, was the first actual procedure to be worked out. The reader familiar with the game of solitaire will recognize the basic similarity.

Suppose that it is desired to form the inverse of a matrix A , of order n . We first form the matrix:

$$E = I - A \quad \dots (7.12.1)$$

and assume, for the present purpose, that the relation:

$$|1 - \lambda_M(A)| = |\lambda_M(E)| < 1 \quad \dots (7.12.2)$$

holds, λ_M being the greatest of the latent roots in the respective matrices. Now subject to equation 7.12.2 it can be shown⁽¹⁹⁾ that:

$$A^{-1} = (I - E)^{-1} = I + E + E^2 + \dots + E^k + \dots = \sum_{k=0}^{\infty} E^k \quad \dots (7.12.3)$$

whence, formally,

$$(A^{-1})_{ij} = \sum_{k=0}^{\infty} (E^k)_{ij} \quad \dots (7.12.4)$$

where $()_{ij}$ represents the element of the i th row and j th column of the matrix $()$.

We now consider a 'random walk' of the following type. Let P_1, P_2, \dots, P_n be a set of (n) points, then we start from any point P_i and jump from point to point in such a way that the probability of a direct move from P_r to P_s is p_{rs} . At any point P_r there is a probability $p_r = 1 - \sum_{s=1}^n p_{rs}$ that the walk ends there. Now, at each transition $P_r \rightarrow P_s$, there is a transition value v_{rs} defined by:

$$v_{rs} \cdot p_{rs} = e_{rs} \quad \dots (7.12.5)$$

where e_{rs} is the corresponding element of the matrix E . If we define the value of the 'walk' to be:

$$W_{ij} = \begin{cases} 0 & \text{if the walk ends at any point } P_k \text{ where } k \neq j \\ v_{i_1 i_1} \cdot v_{i_1 i_2} \cdot \dots \cdot v_{i_m k} \cdot p_j^{-1} & \text{when the walk ends at } P_j \end{cases}$$

where the initial value is unity and the points of the walk are $P_{i_0}, P_{i_1}, \dots, P_{i_m}, P_j$, then the expected value can be shown^(20,21) to be $(A^{-1})_{ij}$; that is the (ij) th element of the required inverse of (A) . This follows from the relation:

(Expectation) $_{ij}$

$$= \delta_{ij} = \sum_{k=1}^{\infty} \sum_{i_1=1}^n \sum_{i_2=1}^n \dots \sum_{i_{k-1}=1}^n e_{i_1 i_1} \cdot e_{i_1 i_2} \cdot \dots \cdot e_{i_{k-1} j} \quad \dots (\text{from 7.12.5})$$

$$= I_{ij} + \sum_{k=1}^{\infty} (E^k)_{ij} = \sum_{k=0}^{\infty} (E^k)_{ij} = (A^{-1})_{ij} \quad \dots (\text{from 7.12.4})$$

where δ_{ij} is the Kronecker delta function, and arises from the fact that when $i = j$ there is a probability p_i that the walk will *at once* stop at P_i , and that the transition value is, in this case, p_j^{-1} .

REFERENCES

- (1) PAIGE, L. J. and TAUSKY, O., 'Simultaneous Linear Equations and the Determination of Eigenvalues,' Bureau of Standards Applied Mathematics Series, No. 29, Washington (1953)
- (2) FERRAR, W. L., 'Algebra, a Textbook of Determinants, Matrices and Algebraic Forms,' Oxford University Press, London (1941)
- (3) BARGMANN, V., MURRAY, F. J. and VON NEUMANN, J., 'Solution of Linear Systems of High Order,' Princeton University Press (1946)
- (4) TODD, J., 'Simultaneous Linear Equations and the Determination of Eigenvalues,' p. 113. *A.M.S.* 29, Washington (1953)
- (5) FERRAR, W. L., 'Algebra,' Theorem 49, p. 153. Oxford University Press, London (1941)
- (6) MORRIS, J., *Phil. Mag.* (7) 37 (1946) 106
- (7) RICHARDSON, L. F., *Phil. Trans. Roy. Soc.*, 242 (1950) 439
- (8) MORRIS, J. and HEAD, J. W., *Phil. Mag.*, 35 (1944) 735
- (9) MORRIS, J., 'The Escalator Method in Engineering Vibration Problems,' Wiley, New York (1947)
- (10) VON NEUMANN, J. and GOLDSTINE, H. H., *Private commun.* (1950)
- (11) FERRAR, W. L., 'Algebra, a Textbook of Determinants, Matrices and Algebraic Forms,' Oxford University Press, London (1941)
- (12) LANCZOS, C., *J. Res. natn. Bur. Stand.*, 45 (1950) 255
- (13) GIVENS, W., *Oak Ridge National Lab. Rept.*, ORNL-1574 (1954)
- (14) HOUSEHOLDER, A. S. and BAUER, F. L., *Num. Math.*, 1 (1959) 29
- (15) WILKINSON, J. H., 'The Algebraic Eigenvalue Problem,' Oxford University Press, London (1965)
- (16) WILKINSON, J. H., *Comput. J.* 3 (1960) 23
- (17) WILKINSON, J. H., *Comput. J.*, 1 (1958) 90
- (18) VON NEUMANN, J. and MORGENTHAU, O., 'The Theory of Games and Economic Behavior,' Princeton University Press (1947)
- (19) COURANT, R. and HILBERT, D., 'Methods of Mathematical Physics,' vol. 1, Interscience, New York (1953)
- (20) FORSYTH, G. E. and LEIBLER, R. A., *Math. Tab., Wash.*, IV (1950) 127
- (21) WASOW, W. R., *ibid.*, VI (1952) 78

PARTIAL DIFFERENTIAL EQUATIONS

8.1 DEFINITIONS AND SCOPE

IN THIS chapter we shall be concerned with the numerical solution of certain simple types of partial differential equation. In the past, it is true to say that only the linear types have received much attention, chiefly because of the technical difficulties of the calculations involved. The present era of high speed automatic digital computers is doing much to remedy this, but the storage capacity of these machines is still too limited to make possible an attack on the more complex types of equation.

The partial differential equations of mathematical physics often fall into the form:

$$A \frac{\partial^2 W}{\partial x^2} + 2H \frac{\partial^2 W}{\partial x \partial y} + B \frac{\partial^2 W}{\partial y^2} + 2G \frac{\partial^2 W}{\partial x \partial z} + 2F \frac{\partial^2 W}{\partial y \partial z} + C \frac{\partial^2 W}{\partial z^2} = I \quad \dots (8.1.1)$$

in which A, B, C, F, G, H, I may be constant (including zero) or functions of (x, y, z) , $\left(\frac{\partial}{\partial x}, \frac{\partial}{\partial y}, \frac{\partial}{\partial z}\right)$. Thus, when

$$A = B = C = 1, H = G = F = I = 0$$

we obtain the Laplace equation:

$$\frac{\partial^2 W}{\partial x^2} + \frac{\partial^2 W}{\partial y^2} + \frac{\partial^2 W}{\partial z^2} = 0.$$

If $I = \text{const.}$ 8.1.1 becomes the Poisson equation, and if W is a function of x and y only, these equations reduce to the well-known two dimensional forms. In two dimensions, if

$$H^2 - AB > 0$$

the equation is said to be *hyperbolic*, if

$$H^2 - AB = 0$$

parabolic, and if

$$H^2 - AB < 0$$

elliptic. These terms are borrowed from the analogous quadratic forms, and are related to convenient subdivisions into which the

PARABOLIC AND HYPERBOLIC EQUATIONS IN TWO VARIABLES

solutions naturally fall. The methods of solution resemble strongly those which have been discussed for ordinary differential equations, especially in their relation to the boundary conditions which the solution has to satisfy. We shall examine, first, simple equations of parabolic and hyperbolic type which are similar in the nature of their solutions to one-point boundary condition equations in the simpler case of ordinary differential equations.

8.2 PARABOLIC AND HYPERBOLIC EQUATIONS IN TWO VARIABLES

The classical example of a parabolic equation represents the one-dimensional flow of heat in a conducting wire, or the diffusion of a liquid or gas along a porous tube, and a study of this equation, from the finite difference viewpoint, lead Courant, Friedrichs and Lewy to their fundamental theorem on the relationship between the space and time intervals. The usual form of the diffusion equation is:

$$c^2 \frac{\partial^2 W}{\partial x^2} = \frac{\partial W}{\partial t} \quad \dots (8.2.1)$$

with initial condition $W = f(x)$ at $t = 0$, and *two* boundary conditions of the general type $W + a \frac{\partial W}{\partial x} = b(t)$ at $x = l_1$, and $x = l_2$, say.

We will assume that the equation is to be solved by replacing both of the derivatives by finite difference approximations at intervals (δx) and (δt) respectively. Using the simple relationships:

$$\frac{\partial^2 W}{\partial x^2} \approx \frac{\delta^2 W}{(\delta x)^2} = \frac{W(x + \delta x, t) - 2W(x, t) + W(x - \delta x, t)}{(\delta x)^2} \quad \dots (8.2.2)$$

$$\text{and: } \frac{\partial W}{\partial t} \approx \frac{\Delta W}{(\delta t)} = \frac{W(x, t + \delta t) - W(x, t)}{(\delta t)} \quad \dots (8.2.3)$$

we obtain, by substitution in equation 8.2.1:

$$W(x, t + \delta t) \approx \gamma W(x + \delta x, t) + (1 - 2\gamma) W(x, t) + \gamma W(x - \delta x, t) \quad \dots (8.2.4)$$

$$\text{where } \gamma = c^2 \delta t / (\delta x)^2. \quad \dots (8.2.5)$$

Equation 8.2.4 enables $W(x, n\delta t)$ to be calculated for any value of n by successive build up *via*

$$W[x, (n - 1)\delta t], W[x, (n - 2)\delta t] \dots W(x, 0).$$

If the value of $W(x, t)$ is required at some particular time T we must take:

$$n\delta t = T$$

which may be considered to determine the interval (δt) in terms of the total time T and the number of steps (n) required. Thus:

$$\gamma = c^2 T / n(\delta x)^2.$$

In a like manner (δx) must be an integral sub-multiple of the length of the (x) boundary.

The important contribution of Courant and his co-workers was to show that it is not possible to choose (δx) and (δt) arbitrarily if a stable solution is to be obtained. By considering the difference, $\epsilon(x, t)$, between the solution of the differential equation 8.2.1 and that of the difference equation 8.2.4, it was shown that this error is bounded only if $\gamma \leq \frac{1}{2}$, and that it grows exponentially with t when $\gamma > \frac{1}{2}$. It follows that, when solving an equation of the general type 8.2.1, (δt) and (δx) must be so chosen as to make $c^2 \delta t / (\delta x)^2 \leq \frac{1}{2}$ and that, in consequence, unlimited decrease in the value of (δx) will not lead to improved accuracy unless accompanied by a suitable decrease in (δt) .

Milne has shown that if $W(x, t)$ has continuous partial derivatives of order 6 in x and of order 3 in t , then the difference between the true solution of equation 8.2.1 and the solution of the difference equation 8.2.4 satisfies:

$$|\epsilon(x, t)| < \frac{c^2 T (\delta x)^4}{135} \left(\frac{\partial^6 W}{\partial x^6} \right)_{\max} \dots (8.2.6)$$

where the value of the partial derivative is taken in the (x, t) region covered by the solution, and the optimum value $\frac{1}{6}$ is taken for γ .

An alternative procedure to the above is to replace only one of the partial derivatives by a finite difference approximation and thus obtain either:

$$\begin{aligned} & c^2 \frac{d^2}{dx^2} [W(x, t + \delta t) + W(x, t)] \\ &= \frac{2}{(\delta t)} [W(x, t + \delta t) + W(x, t)] - \frac{4W(x, t)}{(\delta t)} \dots (8.2.7) \end{aligned}$$

which is an ordinary differential equation in x for $W(x, t + \delta t)$ in terms of the known values $W(x, t)$, or:

$$\frac{c^2}{(\delta x)^2} [W(x + \delta x, t) - 2W(x, t) + W(x - \delta x, t)] = \frac{dW(x, t)}{dt} \dots (8.2.8)$$

which is an ordinary differential equation in t .

The methods described in Chapter 6 can be applied to solve these equations which present no special problems. It may be mentioned that the process defined by equation 8.2.7 is often called the 'Hartree-Womersley' method⁽¹⁾.

A typical hyperbolic equation which occurs in practice is the wave equation:

$$c^2 \frac{\partial^2 W}{\partial x^2} = \frac{\partial^2 W}{\partial t^2} \dots (8.2.9)$$

This time the central difference formulae for the derivatives may be applied to both sides with the result:

$$\begin{aligned} W(x, t + \delta t) &= \beta [W(x + \delta x, t) + W(x - \delta x, t) - W(x, t - \delta t)] \\ &+ 2(1 - \beta)W(x, t) \dots (8.2.10) \end{aligned}$$

$$\text{where} \quad \beta = c^2 \left(\frac{\delta t}{\delta x} \right)^2 \dots (8.2.11)$$

Thus, by taking $\beta = c^2 \left(\frac{\delta t}{\delta x} \right)^2 = 1$, equation 8.2.10 reduces to

$$W(x, t + \delta t) = W(x + \delta x, t) + W(x - \delta x, t) - W(x, t - \delta t) \dots (8.2.12)$$

which has exactly the same solution⁽²⁾, $W = f(x + ct) + g(x - ct)$, as equation 8.2.9. It follows that solutions obtained by replacing equation 8.2.9 by 8.2.10–8.2.12 will be exact. This procedure is particularly useful in dealing with the subsequent motion of a vibrating string whose initial form and velocity are prescribed. In this situation the numerical solution, *via* 8.2.10–8.2.12, is considerably easier to compute than the results of the exact analysis using Fourier series.

The methods just described can be applied in largely unmodified form to parabolic and hyperbolic partial equations of more complicated types, with the exception that the derivatives which occur must be replaced by suitable difference approximations. In a less measure the same remark applies to non-linear equations, although here great care is needed to ensure stability. The reader requiring further information is referred to the excellent paper of BLANCH⁽³⁾.

8.3 HIGHER DIFFERENCES AND CHECKING

A possible method of obtaining a more accurate solution, for the same intervals of differencing, might appear to be to use approximations to the partial derivatives involving higher orders of

differences than those considered in section 8.2. This procedure, unfortunately, often leads to the generation of spurious detail in the resulting solution because, in effect, the original equation is being replaced by a difference equation of higher order. A less dangerous procedure is to use higher difference formulae to estimate the error of a result obtained by the use of the simpler methods.

Thus, in place of equation 8.2.4 we may obtain a more exact relationship *via* the central difference formulae in equations 4.2.10, 4.2.11 of Chapter 4.

$$\frac{\partial^2 W}{\partial x^2} = \frac{1}{(\delta x)^2} (\delta_x^2 - \frac{1}{12}\delta_x^4 + \frac{1}{360}\delta_x^6 - \text{etc.}) W(x, t) \quad \dots (8.3.1)$$

$$\frac{\partial W}{\partial t} = \frac{1}{(\delta t)} [(\mu\delta)_t - \frac{1}{6}\delta_t^2(\mu\delta)_t + \frac{1}{360}\delta_t^4(\mu\delta)_t - \frac{1}{1440}\delta_t^6(\mu\delta)_t] W(x, t) \quad \dots (8.3.2)$$

where $\delta_x \delta_t$ indicate that the differences are to be taken for constant x and constant t respectively. It should be noted that, in making the comparison, the optimum value:

$$\gamma = c^2 \delta t / (\delta x)^2 = \frac{1}{6} \quad \dots (8.3.3)$$

should again be used.

8.4 EQUATIONS WITH MORE THAN TWO VARIABLES

The extension of our previous finite difference treatment to dimensions outside the (x, t) plane presents no special difficulties. We shall illustrate the method of approach by considering the two dimensional Laplacian operator:

$$\nabla^2 W = \frac{\partial^2 W}{\partial x^2} + \frac{\partial^2 W}{\partial y^2} \quad \dots (8.4.1)$$

A fairly obvious start is to replace the partial derivatives by their finite difference approximations:

$$\frac{\partial^2 W}{\partial x^2} \approx \delta_x^2 W / (\delta x)^2 = \frac{1}{(\delta x)^2} [W(x + \delta x, y) - 2W(x, y) + W(x - \delta x, y)]$$

$$\frac{\partial^2 W}{\partial y^2} \approx \delta_y^2 W / (\delta y)^2 = \frac{1}{(\delta y)^2} [W(x, y + \delta y) - 2W(x, y) + W(x, y - \delta y)] \quad \dots (8.4.2)$$

whence:

$$\nabla^2 W \approx \frac{1}{(\delta s)^2} [W(x + \delta x, y) + W(x - \delta x, y) + W(x, y + \delta y) + W(x, y - \delta y) - 4W(x, y)] \quad \dots (8.4.3)$$

where we have assumed that a square difference grid of side $\delta s = \delta x = \delta y$ is being used.

Equation 8.4.3 can be pictured more vividly by the point pattern shown in Figure 8.4.1 (a), which is conveniently represented by the matrix Figure 8.4.1 (b). If we remember that the Laplacian, ∇^2 , of

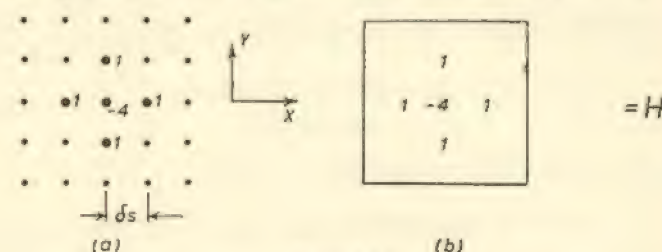


Figure 8.4.1

a function W , at a particular point in space is independent of the particular co-ordinate system chosen, equation 8.4.3 suggests that an alternative expression could be obtained by rotating the axis system through 45° and working in terms of the corresponding points in

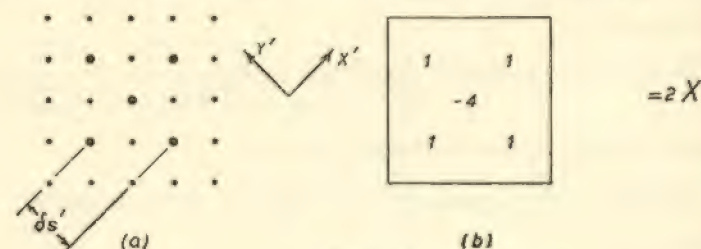


Figure 8.4.2

this system. Inserting the function values in terms of our original co-ordinate system, and noticing that $\delta s' = \delta s \sqrt{2}$, we obtain:

$$\nabla^2 W \approx \frac{1}{2(\delta s)^2} [W(x - \delta x, y + \delta y) + W(x + \delta x, y + \delta y) + W(x - \delta x, y - \delta y) + W(x + \delta x, y - \delta y) - 4W(x, y)] \quad \dots (8.4.4)$$

The preceding arguments have been merely suggestive and we now propose to investigate more closely the nature of the approximations involved; for convenience of writing we shall represent the expressions formed from the function values shown in *Figure 8.4.1 (b)* and *Figure 8.4.2 (b)* respectively, by H and $2X$ ('Horizontal' and 'Cross').

We first introduce the operators E_x and E_y defined by:

$$\left. \begin{aligned} E_x [W(x, y)] &= W(x + \delta x, y) \\ E_y [W(x, y)] &= W(x, y + \delta y) \end{aligned} \right\} \dots (8.4.5)$$

with symbolic associations, corresponding to equation 3.1.11 of Chapter 3,

$$\left. \begin{aligned} E_x &= \exp (\delta x) (\partial / \partial x) \\ E_y &= \exp (\delta y) (\partial / \partial y) \end{aligned} \right\} \dots (8.4.6)$$

From the expression for H we have:

$$H(W) = (E_x + E_x^{-1} + E_y + E_y^{-1} - 4)(W)$$

or, using equation 8.4.6 and the exponential expansion theorem,

$$H(W) = (\delta s)^2 \nabla^2(W) + \frac{1}{12} (\delta s)^4 \nabla_0^4(W) + \frac{1}{360} (\delta s)^6 \nabla_0^6(W) + \dots (8.4.7)$$

where $\nabla_0^n(W) = \frac{\partial^n W}{\partial x^n} + \frac{\partial^n W}{\partial y^n}$. We thus see that a more correct form of 8.4.3 is:

$$\nabla^2 W = \frac{1}{(\delta s)^2} H(W) - O[(\delta s)^2]. \dots (8.4.8)$$

Again, from the expression for $2X$, we have:

$$2X(W) = (E_x^{-1} E_y + E_x E_y + E_x^{-1} E_y^{-1} + E_x E_y^{-1} - 4)(W)$$

which leads to

$$\begin{aligned} 2X(W) &= 2(\delta s)^2 \nabla^2(W) + \frac{1}{6} (\delta s)^4 \nabla_0^4(W) \\ &+ (\delta s)^4 \frac{\partial^4 W}{\partial x^2 \partial y^2} + \frac{(\delta s)^4}{180} \nabla_0^6(W) + \frac{(\delta s)^6}{12} \frac{\partial^4}{\partial x^2 \partial y^2} [\nabla^2(W)] + \dots \end{aligned} \dots (8.4.9)$$

whence:

$$\nabla^2(W) = \frac{1}{2(\delta s)^2} [2X(W)] - O[(\delta s)^2] \dots (8.4.10)$$

and we notice that as is to be expected from the effectively larger interval of differencing involved in calculating $(2X)$, the error term, although $O[(\delta s)^2]$, has a larger coefficient than in equation 8.4.8.

In the special case of the Laplace equation where $\nabla^2(W) = 0$ we have

$$\frac{\partial^4 W}{\partial x^2 \partial y^2} = - \frac{\partial^4 W}{\partial x^4} = - \frac{\partial^4 W}{\partial y^4} = - \frac{1}{2} \nabla_0^4(W)$$

so that equation 8.4.9 becomes:

$$2X(W) = 2(\delta s)^2 \nabla^2(W) - \frac{1}{6} (\delta s)^4 \nabla_0^4(W) + \frac{(\delta s)^6}{180} \nabla_0^6(W)$$

whence, using 8.4.7:

$$(4H + 2X)(W) = 6(\delta s)^2 \nabla^2(W) + \frac{(\delta s)^6}{60} \nabla_0^6(W)$$

$$\text{or} \quad \nabla^2(W) = \frac{1}{6(\delta s)^2} (4H + 2X)W - O[(\delta s)^4]. \dots (8.4.11)$$

The operator $(4H + 2X)$ is easily remembered as:

$$\begin{array}{ccccccc} \cdot & \cdot & \cdot & \cdot & \cdot & & \\ \cdot & \cdot & \cdot & \cdot & \cdot & & \\ \cdot & \cdot & \cdot & \cdot & \cdot & & \\ \cdot & \cdot & \cdot & \cdot & \cdot & & \\ \cdot & \cdot & \cdot & \cdot & \cdot & & \end{array} = \begin{array}{|ccc|} \hline 1 & 4 & 1 \\ \hline 4 & -20 & 4 \\ \hline 1 & 4 & 1 \\ \hline \end{array} = K \dots (8.4.12)$$

the symbol K being suggested by MILNE⁽⁴⁾ who, by introducing the symbolic operator N^2 , defined by:

$$N^2 = 2X - 2H = \begin{array}{|ccc|} \hline 1 & -2 & 1 \\ \hline -2 & 4 & -2 \\ \hline 1 & -2 & 1 \\ \hline \end{array} = (\delta s)^4 \frac{\partial^4}{\partial x^2 \partial y^2} + O[(\delta s)^6] \dots (8.4.13)$$

shows that $\nabla^2 \equiv$

$$\frac{1}{6(\delta s)^2} \left(K - \frac{K^2}{72} + \frac{K^3}{3240} - \frac{KN^2}{180} - \frac{K^4}{120960} + \frac{K^2 N^2}{3780} - \frac{N^4}{504} + \dots \right). \dots (8.4.14)$$

The foregoing analysis enables us to obtain numerical procedures for solving the two-dimensional equivalents of the diffusion and wave equations considered in section 8.2.

Thus, corresponding to equation 8.2.1 we have:

$$\frac{\partial W}{\partial t} = c^2 \nabla^2 W \quad \dots (8.4.15)$$

and obtain at once:

$$W(x, y, t + \delta t) = (\delta t) c^2 \nabla^2 W + W(x, y, t).$$

We now replace $\nabla^2 W$ by a suitable finite difference approximation; if the range of (t) over which the solution is to be continued is *not* large it will be sufficient to use equation 8.4.8. Thus:

$$W(x, y, t + \delta t) = c^2 \frac{\delta t}{(\delta s)^2} HW + W$$

and, if the optimum value $\gamma = c^2 \frac{\delta t}{(\delta s)^2} = \frac{1}{6}$ is chosen, this becomes:

$$W(x, y, t + \delta t) = (\frac{1}{6}H + 1)W \quad \dots (8.4.16)$$

or, symbolically:

$$W(x, y, t + \delta t) = \frac{1}{6} \begin{bmatrix} 1 \\ 1 & 2 & 1 \\ 1 \end{bmatrix} W(x, y, t). \quad \dots (8.4.17)$$

If greater accuracy is required we may use equation 8.4.11 and thus obtain [again with $\gamma = \frac{c^2 \delta t}{(\delta s)^2} = \frac{1}{6}$]:

$$W(x, y, t + \delta t) = (\frac{1}{36}K + 1)W = \frac{1}{36} \begin{bmatrix} 1 & 4 & 1 \\ 4 & 16 & 4 \\ 1 & 4 & 1 \end{bmatrix} W(x, y, t). \quad \dots (8.4.18)$$

The matrix $(K + 36)$ is peculiar in that it can be expressed as:

$$\begin{bmatrix} 1 & 4 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 4 \\ 1 \end{bmatrix}$$

which, as mentioned in a paper by Yowell, quoted by Milne (*loc. cit.*), enables the iterative process defined by equation 8.4.18 to be adapted to calculation on punched card or other fixed scanning sequence computing machines.

The two-dimensional wave equation can be similarly treated, thus:

$$\frac{\partial^2 W}{\partial t^2} = c^2 \nabla^2 W \quad \dots (8.4.19)$$

becomes, on replacing the differentials by differences:

$$W(x, y, t + \delta t) = \frac{c^2 (\delta t)^2}{6 (\delta s)^2} KW(x, y, t) + 2W(x, y, t) - W(x, y, t - \delta t) \quad \dots (8.4.20)$$

Milne has shown that, for convergence, we must take

$$\beta = c^2 (\delta t)^2 / (\delta s)^2 \leq \frac{6}{8}$$

and proposes the value $\beta = \frac{6}{10}$, he thus obtains (actually by using equation 8.4.14 to replace ∇^2)

$$W(x, y, t + \delta t) = \left(\frac{K}{10} - \frac{K^2}{1800} \right) W(x, y, t) + 2W(x, y, t) - W(x, y, t - \delta t) + O[(\delta s)^6] \quad \dots (8.4.21)$$

K^2 is to be obtained by *two* successive applications of K , a procedure which avoids trouble at the boundaries.

Hyperbolic and parabolic partial differential equations are closely analogous to one-point boundary condition ordinary differential equations, in that all the information regarding the solution is known initially, so that the solution can proceed on a step-by-step basis. We have not considered initial conditions of the more complicated types involving first derivatives of the function with respect to space or time co-ordinates, but these can be introduced without any great difficulty, and the diversity of types and situations would make too long an account for the present purpose. In the same way we have considered only hyperbolic and parabolic equations of the most elementary kinds; the method of approach to more complicated examples is still *via* finite difference approximations, but the conditions for convergence are far more difficult to ascertain. Often, in situations representing real problems, a knowledge of the physics of the system will give a clear indication of any marked inaccuracy in the solution; in the absence of such guidance, however, the best that can be done is a check on the solution, as indicated in section 8.3, by means of a more accurate approximation to the equation, using differences of higher order than those used to obtain the solution originally. It cannot be too strongly remarked, however, that such higher difference formulae should *not* be used in the original calculation, since they are equivalent to replacing the original differential

equation by one of higher order, and may thus introduce spurious detail into the solution. A better method is to decrease the size of the interval in the 'open' direction, although here again care must be taken to compensate, if necessary, with a decrease in other intervals. Unfortunately, more general versions of Courant's results are not always available but the known forms may give some idea of the dimensions involved.

8.5 CLASSIFICATION CHARACTERISTICS AND CANONICAL FORM

It was mentioned in section 8.1 that partial differential equations of the type 8.1.1 were divided into classes according to certain criteria; we shall now show how these criteria arise, and how they are related to the canonical forms to which the equation may be reduced.

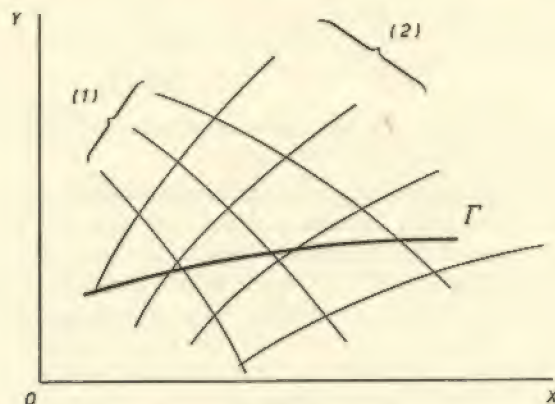


Figure 8.5.1

In the first place we introduce the symbols:

$$\begin{aligned} p &= \frac{\partial W}{\partial x} & q &= \frac{\partial W}{\partial y} \\ r &= \frac{\partial^2 W}{\partial x^2} & s &= \frac{\partial^2 W}{\partial y \partial x} = \frac{\partial^2 W}{\partial x \partial y} & t &= \frac{\partial^2 W}{\partial y^2} \end{aligned}$$

which are common in the theory of surfaces. In terms of these symbols equation 8.1.1 may be written:

$$Ar + 2Hs + Bt = \Phi \quad \dots (8.5.1)$$

where Φ may be a function of the remaining quantities in 8.1.1.

Now the following relationships obtain:

$$dp = \frac{\partial p}{\partial x} dx + \frac{\partial p}{\partial y} dy = r dx + s dy \quad \dots (8.5.2)$$

$$dq = \frac{\partial q}{\partial x} dx + \frac{\partial q}{\partial y} dy = s dx + t dy \quad \dots (8.5.3)$$

Next assume that the values of $\frac{\partial W}{\partial x} (= p)$ and $\frac{\partial W}{\partial y} (= q)$ are given on some curve Γ in the (x, y) plane (Figure 8.5.1). Equations 8.5.1–8.5.3 thus constitute a set of simultaneous equations for the determination of (r, s, t) and will have a determinate solution only if:

$$\Delta = \begin{vmatrix} dx & dy & 0 \\ 0 & dx & dy \\ A & 2H & B \end{vmatrix} \neq 0.$$

The equation $\Delta = 0$ or:

$$A \left(\frac{dy}{dx} \right)^2 - 2H \left(\frac{dy}{dx} \right) + B = 0 \quad \dots (8.5.4)$$

defines two sets of curves. These will be real and distinct only if $H^2 - AB > 0$; only one curve will exist if $H^2 - AB = 0$, and no real curves will exist if $H^2 - AB < 0$; the curves are known as the 'characteristics' of the given equation.

If the two families of characteristics corresponding to the solutions of 8.5.4 are called:

$$\phi(x, y) = \text{const.} \quad \text{and} \quad \psi(x, y) = \text{const.}$$

it may be shown⁽⁵⁾ that the substitution:

$$\xi + i\eta = \phi(x, y), \quad \xi - i\eta = \psi(x, y) \quad \dots (8.5.5)$$

reduces equation 8.1.1 to the canonical *elliptic* form:

$$\frac{\partial^2 W}{\partial \xi^2} + \frac{\partial^2 W}{\partial \eta^2} = X \left(W, \frac{\partial W}{\partial \xi}, \dots, \xi, \dots \right) \quad \dots (8.5.6)$$

whilst the substitution:

$$\xi = \phi(x, y), \quad \eta = \psi(x, y) \quad \dots (8.5.7)$$

gives the canonical *hyperbolic* form:

$$\frac{\partial^2 W}{\partial \xi \partial \eta} = X \left(W, \frac{\partial W}{\partial \xi}, \dots, \xi, \dots \right). \quad \dots (8.5.8)$$

The *parabolic* form is obtained when:

$$\xi = \phi(x, y) = \psi(x, y), \quad \eta = x \quad \dots (8.5.9)$$

and is

$$\frac{\partial^2 W}{\partial \eta^2} = X \left(W, \frac{\partial W}{\partial \xi} \dots, \xi \dots \right). \quad \dots (8.5.10)$$

Characteristics thus specify, in some sense, a set of co-ordinates in which the general equation reduces to the simpler canonical form. It thus becomes appropriate in certain cases to take as a grid for finite differences the characteristics themselves. This technique is likely to be useful only in simple cases in which equation 8.5.4 can be handled by direct analytical methods; it has found some applications in aerodynamics.⁽⁶⁾

8.6 MULTI-POINT BOUNDARY CONDITIONS AND ELLIPTIC PARTIAL DIFFERENTIAL EQUATIONS

Although it may appear curious to include a discussion of *ordinary* differential equations in a chapter on *partial* equations, this position is really the most suitable for those of the former class which have to satisfy multi-point boundary conditions. The essential identity of the principal method of solution for elliptic partial differential equations and that for the ordinary variety will, we hope, be clear after reading this section.

We will start by considering the simplest possible ordinary differential equation of the second order:

$$y'' = 0 \quad (8.6.1)$$

and we will assume that the solution has to satisfy the two-point conditions ($y = y_0, x = 0$) ($y = y_n, x = n\delta x$). Now equation 8.6.1 can be approximated by:

$$\delta^2 y_m = 0 \quad (m = 1, 2 \dots n-1)$$

that is by the set of simultaneous equations:

$$\left. \begin{array}{l} y_0 - 2y_1 + y_2 = 0 \\ y_1 - 2y_2 + y_3 = 0 \\ y_2 - 2y_3 + y_4 = 0 \\ \dots \dots \dots \\ y_{n-3} - 2y_{n-2} + y_{n-1} = 0 \\ y_{n-2} - 2y_{n-1} + y_n = 0 \end{array} \right\} \dots (8.6.2)$$

These equations may be rearranged into the matrix form

$$\begin{bmatrix} 2 & -1 & 0 & 0 & \dots & 0 \\ -1 & 2 & -1 & 0 & \dots & 0 \\ 0 & -1 & 2 & -1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 0 & 0 & 0 & 0 & -1 & 2 \\ 0 & 0 & 0 & 0 & -1 & 2 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ \dots \\ y_{n-2} \\ y_{n-1} \end{bmatrix} = \begin{bmatrix} y_0 \\ 0 \\ 0 \\ \dots \\ 0 \\ y_n \end{bmatrix} \dots (8.6.3)$$

or $A \cdot y = b$ say, whence, if we can find A^{-1} , we can obtain the solution to 8.6.1 as $y = A^{-1} \cdot b$.

Now, in the case of the matrix A of equation 8.6.3 the inverse is easily shown to be:

$$A^{-1} = \begin{bmatrix} (1 - 1/n) & (1 - 2/n) & (1 - 3/n) & (1 - 4/n) & \dots & 2/n & 1/n \\ (1 - 2/n)2 & (1 - 2/n)2 & (1 - 3/n)2 & (1 - 4/n) & \dots & 2.2/n & 2.1/n \\ (1 - 3/n)2 & (1 - 3/n)3 & (1 - 3/n)3 & (1 - 4/n) & \dots & 3.2/n & 3.1/n \\ (1 - 4/n)2 & (1 - 4/n)3 & (1 - 4/n)4 & (1 - 4/n) & \dots & 4.2/n & 4.1/n \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ 1/n & 2/n & 3/n & 4/n & \dots & (n-1).1/n \end{bmatrix} \dots (8.6.4)$$

so that we obtain as our finite difference solution to equation 8.6.1, subject to the given boundary conditions:

$$(y) = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_{n-1} \end{bmatrix} = A^{-1} \cdot \begin{bmatrix} y_0 \\ 0 \\ 0 \\ \dots \\ 0 \\ y_n \end{bmatrix} = \begin{bmatrix} (1 - 1/n)y_0 + y_n/n \\ (1 - 2/n)y_0 + 2y_n/n \\ (1 - 3/n)y_0 + 3y_n/n \\ \dots \dots \dots \\ (1/n)y_0 + (n-1)y_n/n \end{bmatrix} \dots (8.6.5)$$

Since the analytical solution of equation 8.6.1, subject to the boundary conditions, is

$$y = \frac{(y_n - y_0)}{n\delta x} x + y_0$$

we see that the solution given by 8.6.5 is, in this case, *exact*.

The same method can be adopted for the more general equation:

$$y'' = f(x)$$

with similar boundary conditions. In this case the set of difference equations becomes:

$$\delta^2 y_m = f(m\delta x) \cdot (\delta x)^2 \quad (m = 1, 2, \dots, n-1)$$

and the solution may be written:

$$(y) = \begin{bmatrix} y_1 \\ y_2 \\ - \\ - \\ y_{n-1} \end{bmatrix} = A^{-1} \cdot \begin{bmatrix} y_0 - f(\delta x)(\delta x)^2 \\ -f(2\delta x)(\delta x)^2 \\ -f(3\delta x)(\delta x)^2 \\ - \\ y_n - f[(n-1)\delta x](\delta x)^2 \end{bmatrix} \quad \dots (8.6.6)$$

where A^{-1} is the matrix defined by equation 8.6.4.

These simple examples will give the idea behind a general method for solving multi-point boundary condition differential equations. The solution proceeds in three parts:

- (1) Represent the given equation by a finite difference approximation of the same order.
- (2) Set up the system of simultaneous equations defined by the finite difference approximations at each point ($m\delta x$) [$m = 1 \dots (n-1)$] at which a solution is required, using the boundary conditions to define y_0 and y_n in the case of second order equations; for higher orders the boundary conditions may define y_0, y_1 etc. $\dots y_{n-1}, y_n$.
- (3) Solve the simultaneous equations to give the required solution of the original differential equation.

It is evident that the fortunate accident of A having a known inverse will not, in general, happen, so that the set of simultaneous equations, corresponding to 8.6.2, will have to be solved by one of the methods described in Chapter 7. Much existing work in solving two-point boundary condition differential equations has been carried out by the use of the relaxation method (Chapter 7, sections 7.7, 7.9), and the computational layout for such cases has been well described by SOUTHWELL.⁽⁷⁾ When a high speed automatic digital calculator is available, however, modern practice would favour the use of pivotal condensation (Chapter 7, section 7.3). As a check in the work, a finite difference approximation of higher order may be used. This, however, must not be made the basis of the original solution.

Non-linear equations can be solved by the same method, but since these will lead to non-linear simultaneous equations the work of solving the latter may be prohibitive.

A natural extension of the preceding discussion allows the solution of elliptic partial differential equations. We shall consider first the most common form of these—the Laplace equation in two dimensions:

$$\nabla^2 W = 0 \quad \dots (8.6.7)$$

with $W = W_T$ on some boundary Γ . To start with, it will be assumed that Γ consists only of elements parallel to the co-ordinate axes (x, y).

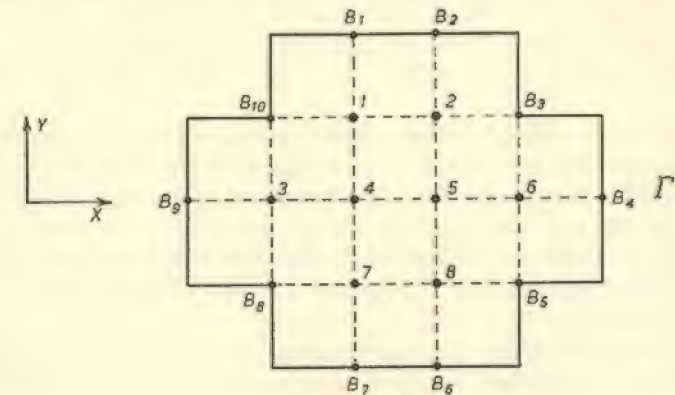


Figure 8.6.1

We divide the interior of Γ over which a solution to equation 8.6.7 is required, by means of a co-ordinate grid of equispaced lines, and call the interior points of this grid, 1, 2, \dots 8 (Figure 8.6.1). At each such point the equation 8.6.7 may be approximated by $H(W) = 0$ (8.4.8) so that, for the region shown in Figure 8.6.1, we obtain the set of relations:

$$H(W_r) = 0 \quad (r = 1 \dots 8)$$

and, if the boundary values are $B_1 \dots B_{10}$ as shown, these become:

$$\left. \begin{aligned} 4W_1 - W_2 - W_4 &= B_1 + B_{10} \\ -W_1 + 4W_2 - W_6 &= B_2 + B_3 \\ &= B_3 + B_9 + B_{10} \\ -W_1 - W_3 + 4W_4 - W_5 - W_7 &= 0 \\ -W_2 - W_3 + 4W_4 - W_5 - W_6 - W_7 - W_8 &= 0 \\ &= B_3 + B_4 + B_5 \\ W_4 + 4W_5 - W_6 - W_7 - W_8 &= B_7 + B_8 \\ -W_5 - W_7 + 4W_8 &= B_8 + B_6 \end{aligned} \right\} \quad \dots (8.6.8)$$

In matrix form this set of equations may be written:

$$G.W = b \quad \dots (8.6.9)$$

where:

$$G = \begin{bmatrix} 4 & -1 & 0 & -1 & 0 & 0 & 0 & 0 \\ -1 & 4 & 0 & 0 & -1 & 0 & 0 & 0 \\ 0 & 0 & 4 & -1 & 0 & 0 & 0 & 0 \\ -1 & 0 & -1 & 4 & -1 & 0 & -1 & 0 \\ 0 & -1 & 0 & -1 & 4 & -1 & 0 & -1 \\ 0 & 0 & 0 & 0 & -1 & 4 & 0 & 0 \\ 0 & 0 & 0 & -1 & 0 & 0 & 4 & -1 \\ 0 & 0 & 0 & 0 & -1 & 0 & -1 & 4 \end{bmatrix} \quad \dots (8.6.10)$$

and b is the column vector whose components are the right-hand elements of equations 8.6.8. It is evident, from the form of equations 8.6.8, 8.6.9, that the matrix G is independent of the boundary values themselves, and also that G is always symmetrical. Further investigation^(a) shows that G is non-singular, and has latent roots which lie in the range $0 < \lambda_i < 8$, these roots being symmetrically located with respect to $\lambda = 4$.

The solution of the Laplace equation 8.6.7 subject to the given boundary conditions is therefore approximated by the solution to 8.6.9:

$$W = G^{-1}.b \quad \dots (8.6.11)$$

so that the methods of Chapter 7, and in particular the relaxation method, can be used to solve this problem. If a solution to the Laplace equation is required for the same boundary *shape* but a number of different boundary *values*, it may be more expeditious to evaluate G^{-1} and then to use equation 8.6.11 rather than to solve the given equations by relaxation each time.

Exactly the same method can be applied to the solution of Poisson's equation:

$$\nabla^2 W = \rho(x, y) \quad \dots (8.6.12)$$

except that the individual difference equations will now become:

$$H(W_r) = (\delta s)^2 \rho_r$$

where ρ_r indicates that the value of ρ is to be taken at the lattice point corresponding to W_r .

Thus, in the example of Figure 8.6.1, the resulting equations are:

$$G.W = \begin{bmatrix} B_1 + B_{10} + (\delta s)^2 \rho_1 \\ B_2 + B_3 + (\delta s)^2 \rho_2 \\ B_8 + B_9 + B_{10} + (\delta s)^2 \rho_3 \\ \text{etc.} \\ B_5 + B_6 + (\delta s)^2 \rho_8 \end{bmatrix} \quad \dots (8.6.13)$$

where G is the matrix of equation 8.6.10.

We now consider the question of lattice dimensions and boundaries which are not straight lines. In the first place, it is clear that (δs)

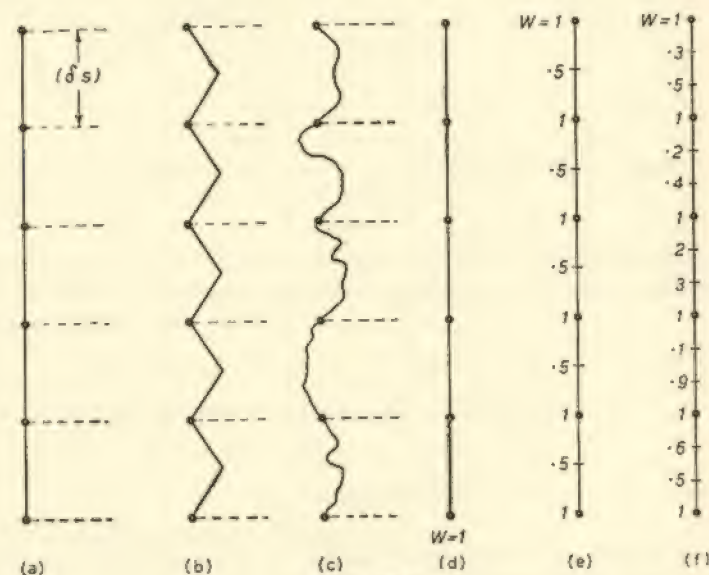


Figure 8.6.2

must be chosen so that variations in boundary shape and value *can* make themselves felt on the solution. Thus, with the mesh (δs) shown in Figure 8.6.2, boundaries (a), (b) and (c) and boundary values (d), (e) and (f) will lead to identical solutions of the finite difference equations approximating 8.6.7 or 8.6.12 although it is quite evident that the true solutions are by no means identical. It is clear, therefore, that the interval (δs) must be so chosen as to allow adequate representation of the detail of the solution as revealed by the known boundary values.

On the other hand, physical considerations, for example, of the related problems of charged conductors, suggest that far from the

boundaries much of this detail will be lost, so that an overall decrease in mesh size would be wasteful in computing time. To overcome this difficulty it is usual to take a small mesh size near to sharp discontinuities and to increase this towards the body of the region of solution. This is illustrated in the enlarged view of the region $B_1 B_{10}$, B_9 of Figure 8.6.1 given in Figure 8.6.3 where it is seen that a mesh of different size has been adopted in the region of the discontinuities.

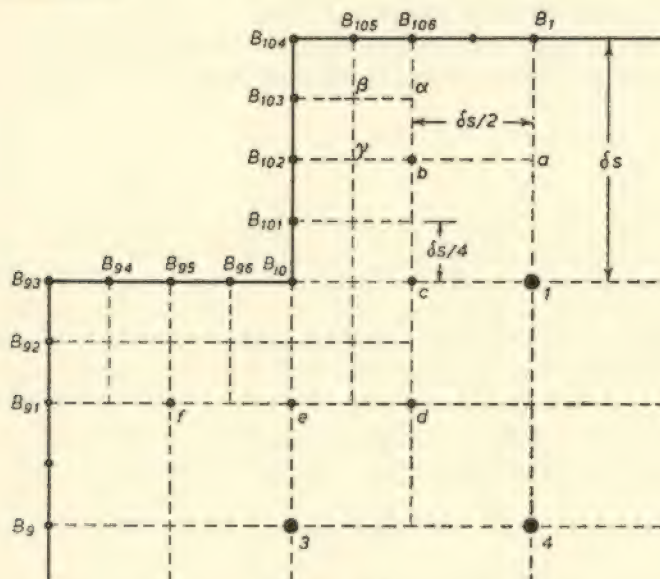


Figure 8.6.3

The joining up of such small meshes to the main one is carried out in steps; for example, in Figure 8.6.3 we should:

- (1) Evaluate the solution using only the coarse mesh shown in Figure 8.6.1.
- (2) Use linear interpolation from (1) and the boundary values B_1 , B_{10} , B_{104} to find approximate values of W at (a) and (b). Thus $W_{(a)} \approx \frac{1}{2}(W_1 + W_{B1})$, $W_{(b)} \approx \frac{1}{2}(W_1 + W_{B1} + W_{B10} + W_{B104})$. Similarly for $W_{(c)} \dots W_{(f)}$.
- (3) Apply relaxation or other techniques to refine these approximations, noticing that alterations in $W_{(a)} \dots W_{(f)}$ will produce changes in the main net $W_1 - W_8$ so that this will require further treatment at the same time.

- (4) When adequately small residuals have been obtained for $W_{(a)} \dots W_{(f)}$, $W_1 \dots W_8$, add the next finer mesh (a) (b) (c) etc. and proceed as before.

As the intervals decrease near the boundary it should be found that less and less alteration is required in the main structure.

The above types of procedure are well adapted to straight boundaries of the sort illustrated, but with curved boundaries, and

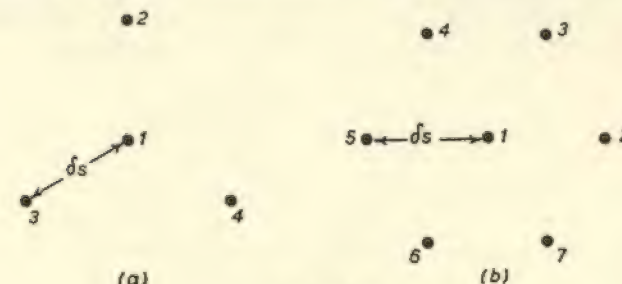


Figure 8.6.4

for straight boundaries which suggest hexagonal or related structures it is sometimes more convenient to adopt one of the finite difference approximations⁽⁹⁾:

$$3W_1 - W_2 - W_3 - W_4 = -\frac{3}{4}(\delta s)^2 \rho_1 \quad \dots (8.6.14)$$

for the triangular mesh shown in Figure 8.6.4 (a), or

$$6W_1 - W_2 - W_3 - W_4 - W_5 - W_6 - W_7 = -\frac{3}{2}(\delta s)^2 [\rho_1 + \frac{1}{16}(\delta s)^2 \nabla^2 \rho_1] \quad \dots (8.6.15)$$

for the hexagonal mesh of Figure 8.6.4 (b).

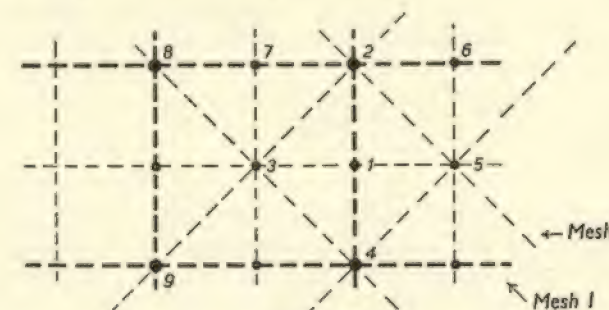


Figure 8.6.5

We may remark, in passing, that it has been recommended that, in decreasing the mesh size, a two stage process of the type shown in Figure 8.6.5 should be adopted:

First the large grid, one element of which is (2, 4, 8, 9), should be refined. Then the diagonal grid (2, 3, 4, 5) should be adopted, the value of W_3 being taken as $\frac{1}{4}(W_2 + W_4 + W_8 + W_9)$. After refinement of this the small grid (1, 2, 7, 3) may be taken with the value W_1 estimated as $\frac{1}{4}(W_2 + W_3 + W_4 + W_5)$. In our experience this cautious approach is not usually justified, and the direct process of halving the mesh edge at each stage is to be preferred.

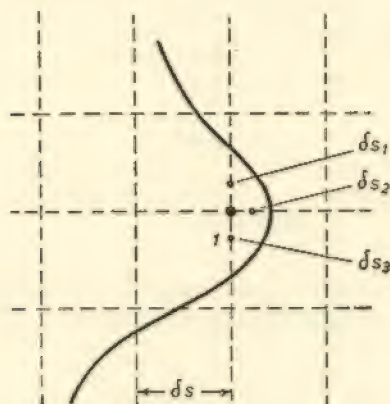


Figure 8.6.6

The general treatment of curved boundaries is difficult; formulae can be constructed to introduce the boundary values in a general situation such as that shown in Figure 8.6.6 for which

$$\nabla^2 W = \frac{2}{\delta s_1 + \delta s_2} \left[\frac{W(x + \delta s_2, y) - W(x, y)}{\delta s_2} + \frac{W(x - \delta s_1, y) - W(x, y)}{\delta s_1} \right] + \frac{2}{\delta s_1 + \delta s_3} \left[\frac{W(x, y + \delta s_1) - W(x, y)}{\delta s_1} + \frac{W(x, y - \delta s_3) - W(x, y)}{\delta s_3} \right] - O(\delta s) \dots (8.6.16)$$

and from which it is seen that the error is considerably greater than is the case with the symmetrical formula 8.4.8. Asymmetrical formulae can be constructed⁽¹⁰⁾ to have an error $O[(\delta s)^2]$ but they are very complicated and we prefer to adopt a locally finer grid in the boundary region and use equation 8.6.16.

8.7 PRACTICAL ASPECTS OF THE RELAXATION METHOD

It seems likely that the large-scale solution of the simultaneous difference equations which result from attempts to solve elliptic partial differential equations will soon be a process which is exclusively performed on high speed computing machines. The best method of solution in this case is one in which a standard procedure is followed in all examples, as, for instance, in the elimination method of section 7.3 or the steepest descent methods of sections 7.7, 7.9.

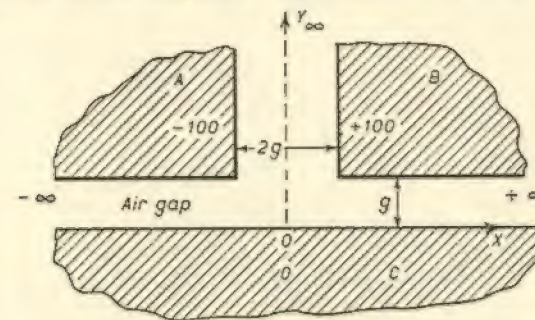


Figure 8.7.1

Occasionally, however, it will always be necessary to solve given problems by hand and in this event the relaxation method seems to offer the simplest procedure and, at the same time, the greatest scope for intuition and experience on the part of the human computer. We therefore propose in this section, to mention briefly, one or two practical aspects of relaxation technique as applied to the solution of equations involving the Laplacian operator.

To illustrate the method we shall examine the two dimensional potential problem represented by the system shown in Figure 8.7.1.

In one physical interpretation of this situation A and B represent the pole faces of a magnetic recording head and C represents the recording medium; all are assumed to define equipotentials.

In obtaining a numerical solution of the Laplace equation

$$\nabla^2 V = 0$$

in the air gap region of Figure 8.7.1 a difficulty arises at once from the fact that the air region extends to infinity in the directions shown. This is overcome in practice by assuming that at some finite distance the solution becomes indistinguishable from the solution at infinity. For the system shown this is the linear function:

$$V_\infty = 100y/g$$

at infinity in a direction parallel to C , and the similar function:

$$V_{\infty} = 100x/g$$

at infinity in a direction parallel to A and B , the y axis being taken to lie midway between A and B .

For the present purpose we will assume that we are interested in a solution of accuracy 1 unit of potential, and experience suggests that the field will be sensibly that at infinity when $OX = 5g$ and similarly for OY . We thus replace Figure 8.7.1 by the system shown in Figure 8.7.2.

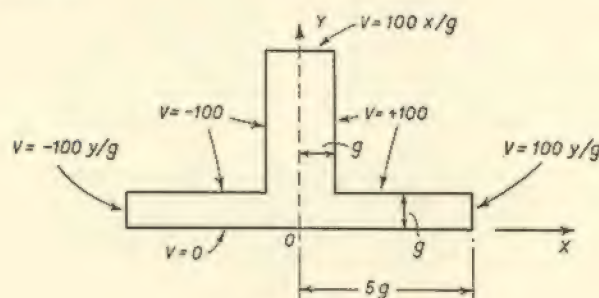


Figure 8.7.2

Next we notice that the distribution will be anti-symmetrical about OY so that OY is at potential zero and Figure 8.7.2 can be reduced to Figure 8.7.3.

Finally we observe that the solution will be symmetrical about OGZ (making an angle of 45° with OX) so that our computations can be confined to the region $OGLX$.

The next consideration is the choice of a finite difference grid upon which to carry out the calculations. Normally, for the accuracy required, we would take a mesh of side $\delta s = g/10$ but for the purpose of this example we shall start with a very coarse mesh ($g/2$) and decrease this progressively to illustrate the technique.

We thus have the system shown in Figure 8.7.4.

The mesh points concerned are $(a \dots i)$ and we might start by taking the initial values of V at these to be zero. If the residuals are now calculated [taking them to be $H(V)$ where H is defined by Figure 8.4.1] we find $H(V_a) = 0$, $H(V_b) \dots H(V_h) = 100$, and $H(V_i) = 150$. Now a situation of this kind in which a number of adjacent mesh points have large residuals of the same sign is usually treated by a technique known as 'block' or 'group' relaxation, a

term which means simply that *all* of the values of V are changed simultaneously and in the same direction until the residuals are reduced as far as possible.

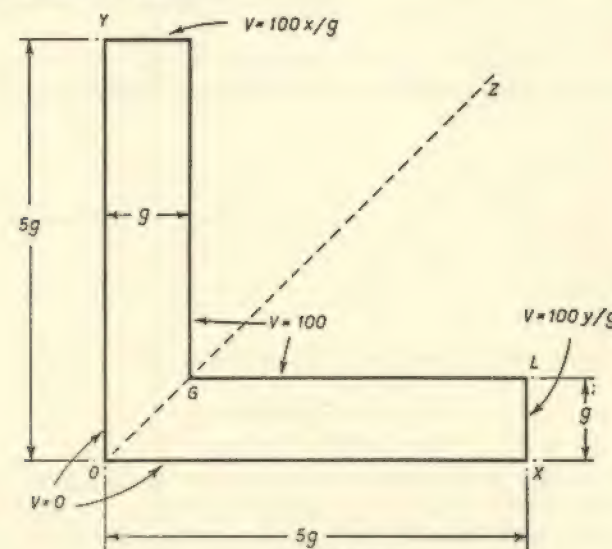


Figure 8.7.3

SOUTHWELL⁽¹¹⁾ has given a physical picture which is often helpful in seeing what sort of movement is required in block relaxation. He shows that the solution is equivalent to finding the equilibrium position of the points $a \dots i$ if they were connected together, and to the boundaries, by means of light elastic strings. The boundaries are assumed to be at levels (perpendicular to the plane of the paper) proportional to their potentials.

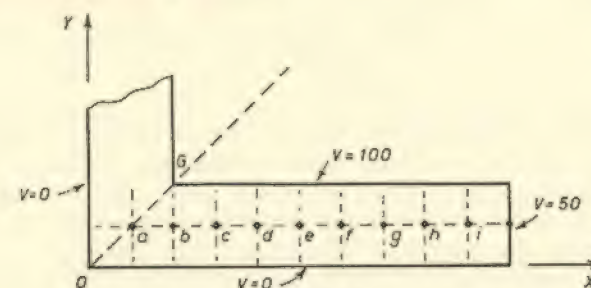


Figure 8.7.4

Thus, for *Figure 8.7.4*, a more plausible starting point would be to take the 'levels' of $a \dots i$ at height 50, i.e. midway between those of the two boundaries. If this is done we find $H(V_a) = -150$, $H(V_b) \dots H(V_i) = 0$ which is considerably more satisfactory.

The first relaxation net is now as shown in *Figure 8.7.5* which also gives the inscriptions which are inserted as the relaxation proceeds.

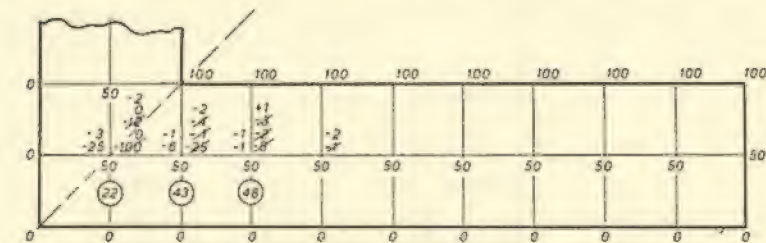


Figure 8.7.5. Relaxation No. 1

It is our practice to write the initial value under each lattice point and then inscribe residuals in a vertical column to the right of the point; corrections, if any, are written on the left. One point of procedure may be mentioned: *remember to correct residuals at adjacent points after altering any lattice-value*, neglect of this obvious warning has often led to error. Numbers are always chosen so that only integers

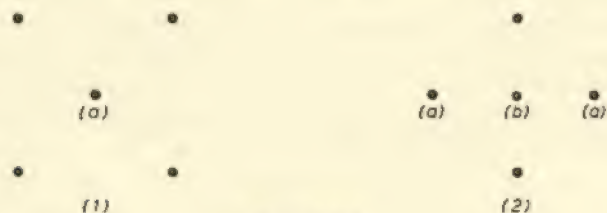


Figure 8.7.6

appear on the computation sheet, this means that it is not worth reducing residuals to below 2 units in magnitude. At this point the corrections are added into the initial values and the residuals recalculated to ensure removal of casual errors, then the process is restarted with the residuals multiplied by 10 or 100.

For the present example we do not yet perform this scaling up of residuals, we first divide the mesh interval by 2, obtaining

Figure 8.7.7, and then refine this as far as possible. The procedure for obtaining initial values at the new points is shown in the figure.

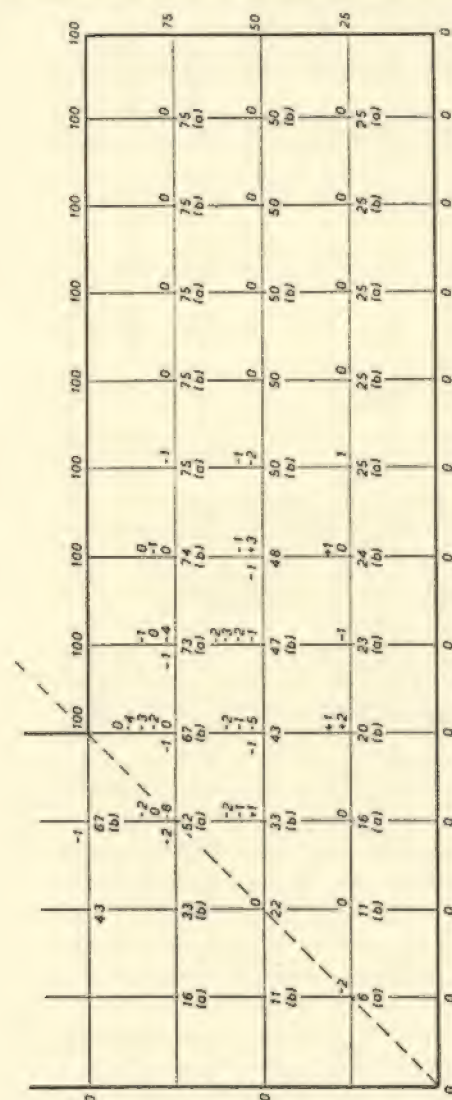


Figure 8.7.7. Relaxation No. 2

first those values marked (a) are derived by taking $\frac{1}{4} \times$ sum of the adjacent net points (and boundary values if necessary) obtained in the first approximation as shown in *Figure 8.7.6 (1)*. When all

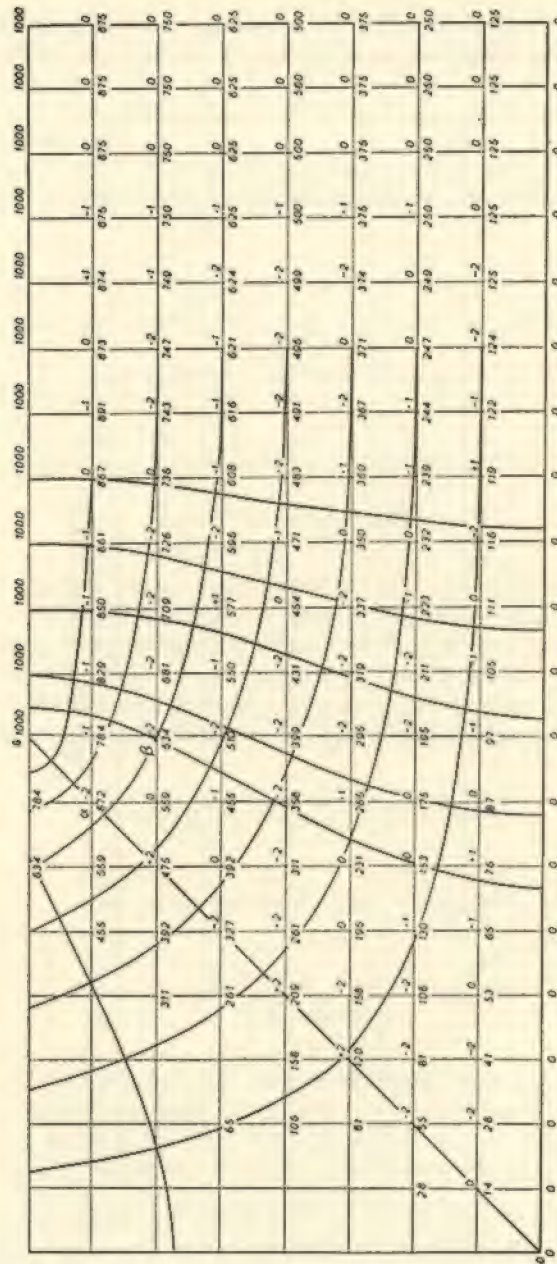


Figure 8.7.8

(a)-type points have been evaluated these are used, with the original values, to obtain intermediate points (b) shown in Figure 8.7.6 (2). The advantage of this procedure is that the (b) points have a maximum residual of ± 2 when we work to the nearest whole number. When the new approximation is complete it is 'relaxed' as before to reduce all residuals to 2 or less.

At this point the interval is again halved and the values corrected, by relaxation, to have residuals less than or equal to two. The results of this process are shown in Figure 8.7.8.

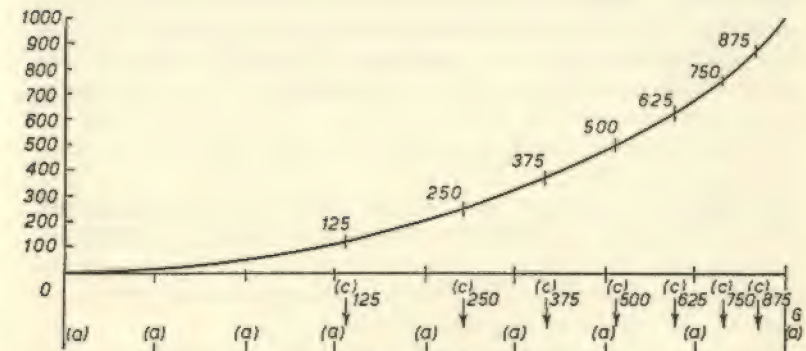


Figure 8.7.9

For the present purpose this approximation is sufficiently exact, but if greater precision were needed the mesh size could again be halved and the process repeated. In Figure 8.7.8 a set of equipotentials have been drawn and, as these are frequently required, we may remark that they can be located with sufficient accuracy by drawing a simple graph of the potential values along each net line and finding the intersections of this curve with ordinates whose values are those of the potentials required. An example of such a graph, for the symmetry axis OG , is shown in Figure 8.7.9. In this example, where equipotentials intersect OG normally, it is a good plan to find the intersections with this line although, in general, such an oblique line would not be needed. As a practical point in technique it should be noticed that the graph for an oblique line such as OG , can be most conveniently drawn by placing the ruled paper with its edge along OG and marking the lattice point intersections as shown in (a) . . . (a). The contour crossings are also marked on the graph at (c) . . . (c) and can be inserted on the grid by replacing the completed graph on OG .

8.8. CHECKING

An advantage of the relaxation method is the fact that errors are automatically corrected although, of course, they hold up the convergence of the process. It is particularly necessary to recompute the residuals *ab initio* periodically since it is rather easy to make mistakes in writing down corrections to adjacent residuals after altering any point.

The estimation of the overall accuracy of the solution is more difficult, it is *not* sufficient to assume that when residuals have been reduced to ± 2 units this is the overall accuracy of the solution, since no account has been taken of the intrinsic errors due to the finite difference approximation. A good guide is available when several grids, of different mesh sizes (δs), have been used in the calculation; generally the error of the final solution at any point is less than the difference between the values obtained, at that point, on the final and penultimate nets.

An analytic estimate of error can be obtained by computing the value of the operator K (8.4.12) at each point, the value $-K/20$ will then give the correction to be applied if K (instead of H) were the relaxation operator, since K is a better approximation to ∇^2 than H (8.4.11), this correction is a measure of the error.

Thus, in our previous example (Figure 8.7.8), the application of K at the points α and β leads to values $R_\alpha = +47$, $R_\beta = -38$ suggesting corrections $+2$, -2 . Since these are in accord with our previous estimates we may take the values given in Figure 8.7.8 as certainly reliable to 1 per cent as required.

Other measures of the accuracy of a solution can be derived. For example, if we assume that errors are of two kinds:

- (1) Those due to the approximations involved in solving the finite difference equations.
- (2) Those due to approximating the operator ∇^2 by the operators H or K ,

then it can be shown⁽¹²⁾ that, if the region over which the solution is computed can be enclosed in a circle of radius R , then the maximum error in the approximate solution of the difference equations is:

$$\epsilon_{(1)} \leq |r_m| R^2 / \epsilon (\delta s)^2 \quad \dots (8.8.1)$$

where r_m is the maximum value of the residual at any point, (δs) is the mesh size, and $\epsilon = 4$ for the operator H and $\epsilon = 24$ for the operator K .

Errors due to cause (2) are more difficult to estimate, but the formula

$$\epsilon_{(2)} \leq \frac{R^2}{24(\delta s)^2} |N^2 W|_{\max} \quad \dots (8.8.2)$$

(where N^2 is the operator defined in equation 8.4.13) gives an order of magnitude in the case of operator H whilst for K

$$\epsilon_{(2)} \leq \frac{R^2}{21096(\delta s)^2} |N^4 W|_{\max} \quad \dots (8.8.3)$$

is appropriate. Neither result is valid for boundary points; 8.8.2 holds, in general, for both the Laplace and the Poisson equations, and 8.8.3 for the Laplace equation only. For the special value $\rho(x, y) = \text{const.}$, in equation 8.6.12, 8.8.3 may also be used.

8.9 MONTE CARLO METHODS

Monte Carlo methods can be derived for the solution of linear partial differential equations and have the advantage that, by their aid, the

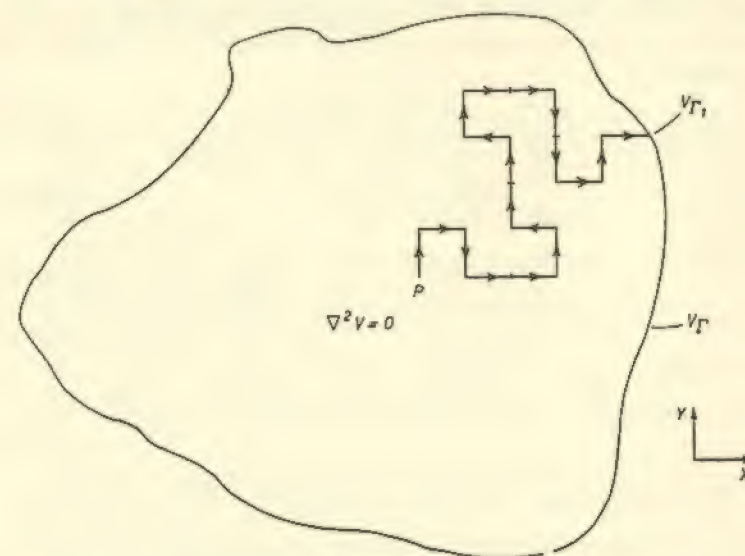


Figure 8.9.1

solution can be obtained at a single point without the need for working over the whole field. We give as an example the solution of the Laplace equation.

Assume that the solution of $\nabla^2 V = 0$ is required at some point P , inside a boundary Γ upon which the values of V are prescribed. Starting at P we make a random walk in which the steps, of length (δs) and in a positive or negative direction, are made parallel to the X or Y axes. The walk continues until it reaches Γ when it is terminated, and the value of V , V_P , say, is recorded. This process is repeated many times (n) and it may be shown ⁽¹³⁾ that:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{r=1}^n V_{P_r} \rightarrow V_P. \quad \dots (8.9.1)$$

The method can be extended to more than two dimensions and also to the case of the Poisson equation; it seems more appropriate to the determination of the values of V at a single point than to its evaluation over the whole interior of Γ , since the number of steps in each walk is of the order $(R/\delta s)^2$ where R is the radius of a circle completely enclosing Γ , and the rate of convergence is proportional only to $1/\sqrt{n}$.

8.10 MORE COMPLICATED PARTIAL DIFFERENTIAL EQUATIONS

The methods outlined in this chapter can be extended, formally, to more complicated equations without difficulty. For example the equations:

$$\frac{\partial w}{\partial t} = c^2 e^x \left(\frac{\partial^2 w}{\partial x^2} + \frac{\partial w}{\partial x} \right) \quad \dots (8.10.1)$$

$$\frac{\partial w}{\partial t} = \frac{\partial^2}{\partial x^2} (w e^{-w}) \quad \dots (8.10.2)$$

$$\frac{\partial w}{\partial t} = c^2 \left(\frac{\partial^2 w}{\partial x^2} + \frac{\partial^2 w}{\partial y^2} + \frac{\partial^2 w}{\partial z^2} \right) \quad \dots (8.10.3)$$

$$\nabla^2 \nabla^2 w = 0 \quad \dots (8.10.4)$$

$$\frac{\partial}{\partial x} \left(\chi \frac{\partial \psi}{\partial x} \right) + \frac{\partial}{\partial y} \left(\chi \frac{\partial \psi}{\partial y} \right) + Z(x, y) = 0 \quad [\chi = \chi(x, y)] \quad \dots (8.10.5)$$

$$\left. \begin{aligned} \frac{\partial v}{\partial t} &= \frac{\partial^2 v}{\partial x^2} + e^{-1/v} & x \geq 0, t > 0 \\ -\frac{\partial v}{\partial x} &= H[v_g(t) - v] & x = 0, t > 0 \\ \frac{\partial v}{\partial x} &\rightarrow 0 & x \rightarrow \infty, t \geq 0 \\ v &= v^{(0)} & x \geq 0, t = 0 \end{aligned} \right\} \quad \dots (8.10.6)$$

$$\left. \begin{aligned} H \frac{\partial u}{\partial x} + K \left(\frac{\partial u}{\partial r} + \frac{\partial v}{\partial x} \right) + L \frac{\partial v}{\partial r} + P &= 0 \\ \frac{\partial v}{\partial x} - \frac{\partial u}{\partial r} &= 0 \\ H = a^2 - u^2, \quad K = -uv, \quad L = a^2 - v^2, \quad P = av^2/r \\ a^2 &= \frac{1}{2}(g-1)(1-u^2-v^2) \end{aligned} \right\} \quad \dots (8.10.7)$$

appear in the literature and methods of approach can be worked out for almost any equation or system of equations which may be suggested. The difficulty is not in evolving systematic means of iterative 'solution', but of ensuring that the results obtained do, in fact, bear some relation to the true answer.

The reader interested in following up the equations just quoted will find them in the references given below:

- 8.10.1-3 MILNE, W. E., 'Numerical Solution of Differential Equations,' Wiley, New York (1953)
- 8.10.4 MILNE, W. E., 'Numerical Solution of Differential Equations,' Wiley, New York (1953)
- 8.10.5 SOUTHWELL, R. V., 'Relaxation Methods in Theoretical Physics,' Oxford (1946)
- 8.10.6 LANDAU, H. G. and HICKS, B. L., *Math. Tab. Wash.*, III (1948) 207
- 8.10.7 LOTKIN, M., *ibid.*, III (1948) 209

REFERENCES

- (1) HARTREE, D. R., and WOMERSLEY, J. R., *Proc. Roy. Soc. A*, 161 (1937), 363
- (2) BOOLE, G., 'Calculus of Finite Differences,' 2nd edn. p. 270. Macmillan, London (1872)
- (3) BLANCH, G., *J. Res. natn. Bur. Stand.*, 50 (1953) 343
- (4) MILNE, W. E., 'Numerical Solution of Differential Equations,' p. 133, Wiley, New York (1953)
- (5) SOMMERFELD, A., 'Partial Differential Equations of Physics,' p. 38 Academic Press, New York (1949)
- (6) SHAPIRO, A. H. and EDELMAN, G. M., *J. Appl. Mech.* A14 (1947) 154
- (7) SOUTHWELL, R. V., 'Relaxation Methods in Engineering Science,' Oxford University Press, London (1940)
- (8) MILNE, W. E., 'Numerical Solution of Differential Equations,' Wiley, New York (1953)
- (9) SOUTHWELL, R. V., 'Relaxation Methods in Theoretical Physics,' Oxford University Press, London (1946)
- (10) MILNE, W. E., 'Numerical Solution of Differential Equations,' p. 150, Wiley, New York (1953)
- (11) SOUTHWELL, R. V., 'Relaxation Methods in Theoretical Physics,' Oxford University Press, London (1946)
- (12) MILNE, W. E., 'Numerical Solution of Differential Equations,' p. 217, Wiley, New York (1953)
- (13) CURTISS, J. H., 'Monte Carlo Methods for the Iteration of Linear Operators', *Nat. Bur. Stand. Rep.* 2365, Washington (1953)

NON-LINEAR ALGEBRAIC EQUATIONS

9.1 PRELIMINARY IDEAS

THE subject of this chapter is the solution of those types of non-linear equation, both simple and simultaneous, which do *not* involve the operations of the differential or integral calculus. The problem may be expressed as:

'Find the solutions of the set of (n) simultaneous equations

$$f_i(x_1, x_2, \dots, x_r, \dots, x_n) = 0 \quad (i = 1 \dots n) \quad \dots(9.1.1)$$

where the function f_i may vary from equation to equation.

If we take a single function, f say, of a single variable, x , the problem is simply:

$$\text{'solve } f(x) = 0\text{'}. \quad \dots(9.1.2)$$

Clearly equation 9.1.1 embraces 9.1.2 and also the case in which the f_i are simple linear combinations of the x_i , the latter case leading to the simultaneous linear equations of Chapter 7.

It will be appropriate to consider, too, those cases in which the equations 9.1.1. form both under- and over-determined sets, a situation which is common in various branches of physical science.

9.2 GRAPHICAL METHODS

The solution of non-linear algebraic equations in a single variable can generally be best performed in two stages. The first may be called a *survey* in which the general characteristics of the system are investigated and in which the roots are located approximately. The second stage may be called the *refinement*, in which the roots are found to the desired degree of accuracy.

It cannot be too strongly urged that even in the simplest problems the survey should be carried out, since the *feel* of the problem, thus obtained, will invariably lead to a better use of the available resources at the refinement stage.

Since the advent of high speed digital calculators the technique of survey has been slightly altered and, in some cases, the mathematical work which had to be carried out as a preliminary to the old methods

GRAPHICAL METHODS

has been thereby eliminated. Nevertheless the availability of the new machines should not be made the excuse for rushing into computation rather than contemplating the problem.

To take a concrete example we may consider the solution of the cubic equation:

$$ax^3 + bx^2 + cx + d = 0. \quad (9.2.1)$$

First we notice that this can be reduced to the form:

$$x^3 + b'x^2 + c'x + d' = 0 \quad \dots(9.2.2)$$

by division through by a , and for a computing machine such a preliminary reduction would almost always be wise. If the survey is to be conducted by hand, however, we should consider the nature of the original coefficients a, b, c, d ; should these be simple whole numbers they may be far easier to manipulate than the, possibly, infinite decimals of equation 9.2.2. We assume that, in the present case, the latter form is appropriate, and consider the next step. On an automatic machine this might consist of a simple evaluation of 9.2.2 for a range of values of x . Considerations of the relative sizes of b', c', d' will limit this to manageable proportions, and may give some clue as to the interval of tabulation required. As a minor point of procedure it may be noted that 9.2.2 should *not* be calculated by forming (on an automatic computer) successively $x^3, b'x^2, c'x$ and adding these to d' , but by means of the sequence:

$$x + b', \quad x(x + b') + c', \quad x[x(x + b') + c'] + d'$$

the reason being that the former process requires four multiplications to form $x^3, b'x^2$ and $c'x$, whereas the latter requires only two.

When an automatic machine is *not* available, however, the above method is not to be recommended. A better approach is to notice that the roots of 9.2.2 are the same as the x co-ordinates of the intersections of

$$\left. \begin{aligned} y &= -c'x - d' \\ y &= x^3 + b'x^2 \end{aligned} \right\}. \quad \dots(9.2.3)$$

The first of these equations represents a straight line, which may be conveniently drawn through the two points: $(0, -d')$ $(-d'/c', 0)$ whilst the second equation is most easily plotted by means of a table of squares and a slide rule—the first to find x^2 and the second to multiply this by $(x + b')$ which is formed mentally.

Figure 9.2.1 (a) (b), illustrates the two methods of approach applied to the simple cubic equation:

$$x^3 - 3x^2 + 4x - 2 = 0.$$

Graph (a) is obtained by a direct plot of:

$$y = x^3 - 3x^2 + 4x - 2$$

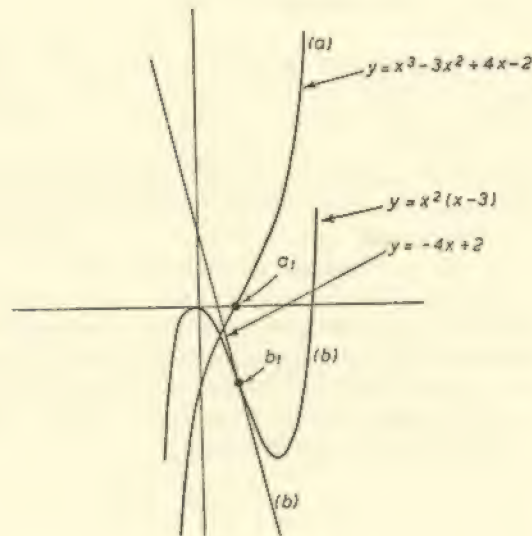


Figure 9.2.1

at interval $x = 0.5$, whilst the pair of graphs (b) give

$$y = -4x + 2$$

$$y = x^2(x - 3)$$

and were plotted from points at the same interval, both reveal the real root at $x = 1$.

The above example is an extremely simple one and does not show the wealth of detail which can sometimes result from the preliminary survey; a more interesting illustration is provided by the equation:

$$ax = \tan x \quad \dots (9.2.4)$$

which occurs in the design of gears. We replace the single equation by:

$$\left. \begin{aligned} y &= ax \\ y &= \tan x \end{aligned} \right\} \quad \dots (9.2.5)$$

and thus obtain the graph of Figure 9.2.2.

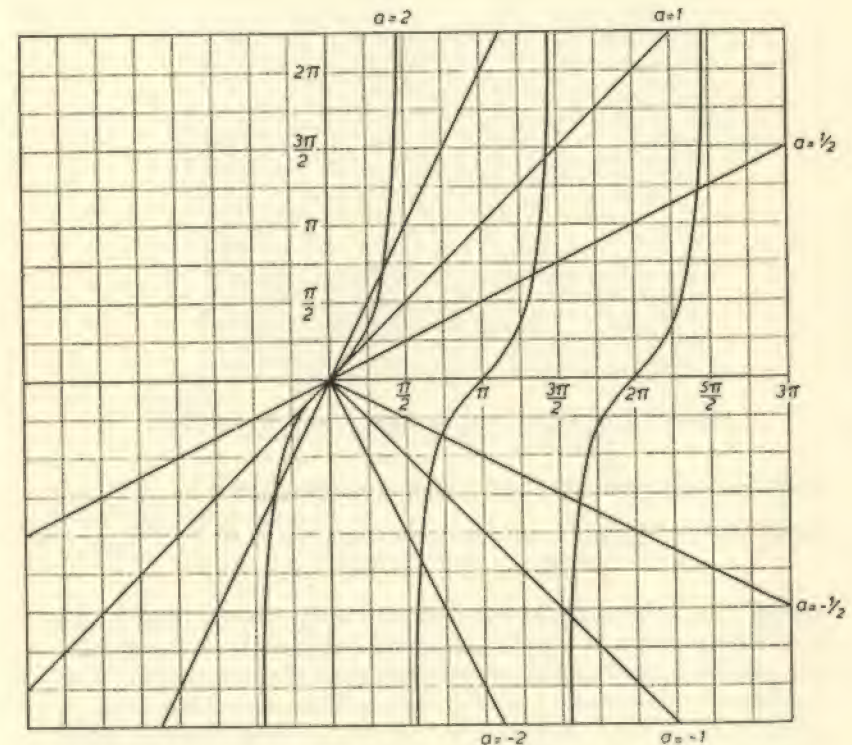


Figure 9.2.2. $\tan x = ax$

The intersections, for various values of a , give the approximate positions of the required roots which can be refined by the processes to be described in section 9.3. In addition to this data, the graphs show that there are no positive real roots (other than $x = 0$) in the range $(0, \pi)$ when $0 < a \leq 1$, and no roots in the range $(0, \pi/2)$ when $a < 0$. Furthermore, it is clear that for large x the roots tend asymptotically to $(n + \frac{1}{2})\pi$ for all values of $|a| > 0$.

An example of a survey applied to a set of simultaneous equations is provided by x-ray crystal structure analysis. It became necessary

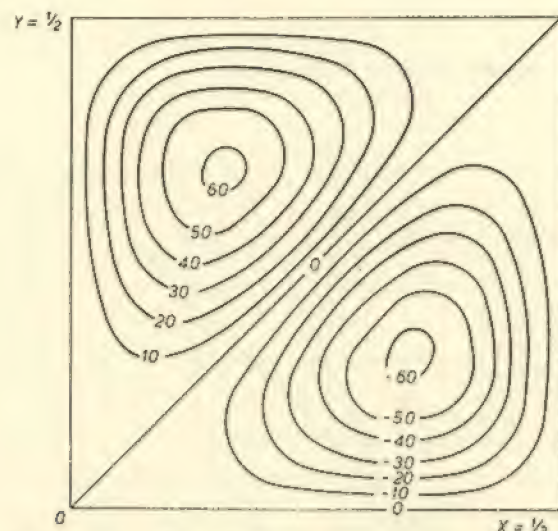


Figure 9.2.3 (a) $40 [\sin (4\pi x) \cdot \sin (2\pi y) - \sin (2\pi x) \cdot \sin (4\pi y)]$

to obtain the solution of the simultaneous equations:

$$\left. \begin{aligned} 40[\sin (4\pi x) \cdot \sin (2\pi y) - \sin (2\pi x) \cdot \sin (4\pi y)] &= -30 \\ 40[\cos (8\pi x) + \cos (8\pi y)] &= +10 \end{aligned} \right\} \dots (9.2.6)$$

and for these the two dimensional contour maps shown in Figure 9.2.3 (a) (b) are plotted. By superimposing the two maps, as in Figure 9.2.4 it is possible to read off the approximate roots of the equations 9.2.6, and also to see exactly how the various alternative pairs of (x, y) values arise, a point which is by no means clear from a direct analytic solution.

For systems of equations which involve more than two variables the graphical survey is generally impossible. In these cases, apart from any information as to the locations of roots which may result from a knowledge of the experimental system in which they arise, there will usually be no alternative but to take some arbitrary starting point and use one of the successive refinement methods to be described in the remainder of this chapter. A word of warning is perhaps appropriate in this connection; it is that refinement methods usually postulate that the trial point is near to the true solution. If

this condition is not satisfied, they may either not converge at all or, more probably, converge to a false solution which is in the nature

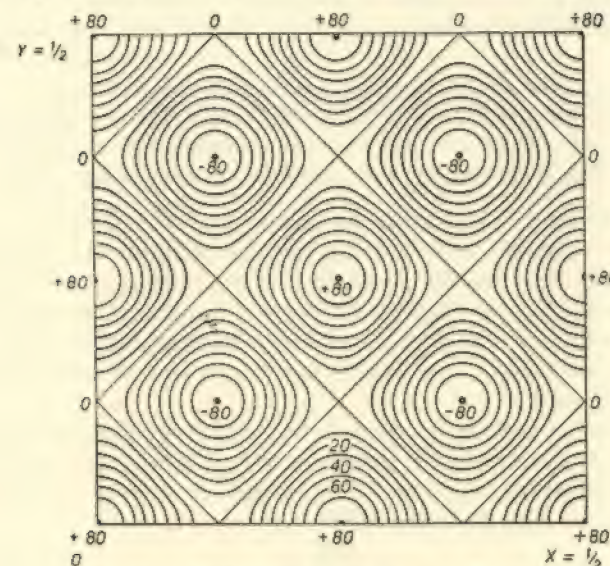


Figure 9.2.3 (b) $40 [\cos (8\pi x) + \cos (8\pi y)]$

of a 'col', or depression, in a mountain range, instead of the true valley which is sought.

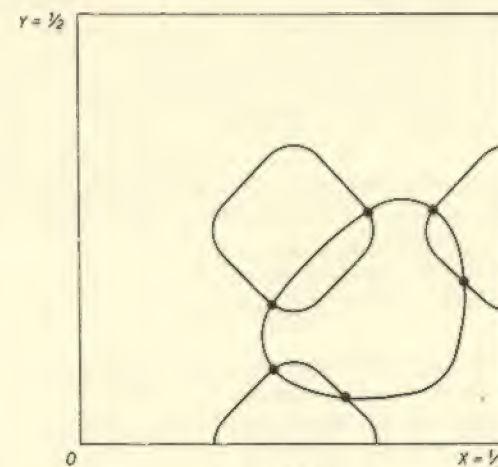


Figure 9.2.4

9.3 ITERATIVE PROCESSES—ONE VARIABLE

We may, as in section 9.1, write the general non-linear equation in one variable as:

$$f(x) = 0 \quad \dots (9.3.1)$$

and the problem which we shall consider in this section is the solution of this form. We may remark that when $f(x)$ is a polynomial, say:

$$f(x) = x^n + a_1x^{n-1} + a_2x^{n-2} \dots + a_n = 0 \quad \dots (9.3.2)$$

classical methods of solution, depending upon the theory of equations, are available. We shall not mention them further except to remark that they are the methods of Horner and Graeffe^(1, 2) and that they find little favour with modern numerical analysts. The analytic solutions of the cubic and diquadratic equations, known as Ferrari's, Tartaglia's and Cardan's methods, are of little utility in numerical work.

An iterative process for the solution of an equation of the type 9.3.1 may be defined by the statement that, if x_0 is an approximation to the solution of equation 9.3.1, it enables a quantity x_1 to be calculated by means of some relation:

$$x_1 = I(x_0) \quad \dots (9.3.3)$$

in such a manner that x_1 is a closer approximation to the required solution than was x_0 .

If x_0 differs from the true root by a small quantity, of order (ϵ) , say, then the iterative process 9.3.3 is said to be n th order if the error in x_1 is of order (ϵ^n) . Most of the best iterative processes are second order; third and higher order processes exist and can always be constructed from those of lower order, but often they involve greater *total* computing labour for a given final accuracy than those of the second order.

As an example of a first order process we may consider the so-called *regula falsi* or rule of false position which is possibly the oldest known iterative procedure.

Suppose it is required to find the root of:

$$y = f(x) = 0$$

corresponding to the intersection R shown in Figure 9.3.1. Any point (x_0, y_0) on the curve is taken and also any other point (x_1, y_1) —preferably on the opposite side of OX from (x_0, y_0) . Effectively the process consists of finding the intersection of the line joining (x_0, y_0) and (x_1, y_1) with OX and then taking this abscissa, x_2 say, to determine a new

point (x_2, y_2) on the curve. The operation sequence is repeated using (x_0, y_0) and (x_2, y_2) , and it is clear that convergence to the root, R say, takes place. Analytically the method is equivalent to a sequence of linear interpolations, and the relation between successive convergents is easily shown to be:

$$x_{n+1} = (x_0y_n - x_ny_0)/(y_n - y_0). \quad \dots (9.3.4)$$

To investigate the relation between the errors of the n th and $(n+1)$ th approximation assume the root to be $x = X$ and let:

$$x_0 = X + \epsilon_0$$

$$x_n = X + \epsilon_n$$

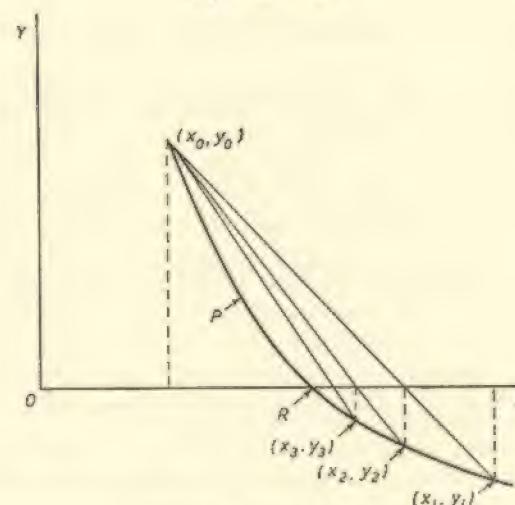


Figure 9.3.1

Then Taylor's theorem gives:

$$y_0 = \epsilon_0 f'(X) + \frac{1}{2} \epsilon_0^2 f''(X) + \dots$$

$$y_n = \epsilon_n f'(X) + \frac{1}{2} \epsilon_n^2 f''(X) + \dots$$

(since, by definition $f(X) = 0$) and substitution in 9.3.4 and simplification leads to:

$$x_{n+1} \approx X + \frac{1}{2} \epsilon_0 \cdot \epsilon_n f''(X) / f'(X)$$

whence

$$\epsilon_{n+1} \approx \frac{1}{2} \epsilon_0 \cdot \epsilon_n f''(X) / f'(X) \quad \dots (9.3.5)$$

thus, if ϵ_0 is such that $|\frac{1}{2} \epsilon_0 f''(X) / f'(X)| < 1$, the process will

converge. Actually convergence may be quite rapid if the initial point is well chosen, but it should be noticed that each iteration of the process will only produce an improvement of the order of one place of decimals.

Another well-known technique for solving polynomial equations in one variable is the method of iteration. The equation to be solved, $f(x) = 0$, is written in the form

$$x = g(x)$$

The iterative process is then simply:

$$x_{n+1} = g(x_n)$$

Clearly, if the process converges $\mathcal{L}t$ x_n tends to the solution required.

To investigate the convergence of the process assume that the approximation, x_n , differs from the desired root α by a small quantity ϵ . Then, since α is a root of $f(x) = 0$,

$$\alpha = g(\alpha)$$

Also $x_{n+1} = g(x_n) = g(\alpha + \epsilon)$

Whence, expanding, by Taylor's theorem, and neglecting terms in ϵ^2 etc.,

$$\begin{aligned} x_{n+1} &= g(\alpha) + \epsilon g'(\alpha) \\ &= \alpha + \epsilon g'(\alpha) \end{aligned}$$

It follows that the process converges if $|g'(\alpha)| < 1$.

The advantage of the iterative method is its simplicity, its disadvantage is that, since α is initially unknown, it must often happen that convergence cannot be determined *a priori*.

A general technique of construction for second order processes is that known as the Newton-Raphson method. Assume that an approximation (x_n) has been calculated which has an error (ϵ_n), then since:

$$f(x_n + \epsilon_n) = 0$$

we have, by Taylor's theorem:

$$f(x_n + \epsilon_n) = 0 = f(x_n) + \epsilon_n f'(x_n) + \frac{\epsilon_n^2}{2} f''(x_n) + \dots$$

so that, if we choose

$$x_{n+1} = x_n - f(x_n)/f'(x_n) \quad \dots (9.3.6)$$

the error in x_{n+1} is of order $(\epsilon_n)^2$.

It should be noticed that there is no *unique* iteration of the second order leading to a given result. Thus, suppose that a process for determining the square root of a quantity b is required, the equation:

$$x^2 - b = 0 \quad \dots (9.3.7)$$

has the required solution $x = \pm \sqrt{b}$.

From equation 9.3.6 the iteration:

$$x_{n+1} = x_n - (x_n^2 - b)/2x_n = \frac{1}{2}(x_n + b/x_n) \quad \dots (9.3.8)$$

is obtained. On the other hand, 9.3.7 may be written:

$$1/x^2 - 1/b = 0$$

and this leads to:

$$x_{n+1} = x_n - (1/x_n^2 - 1/b)/(-2/x_n^3) = x_n[1 + (b - x_n^2)/2b]. \quad \dots (9.3.9)$$

It can be shown that, although equations 9.3.8. and 9.3.9 are both second order processes, equation 9.3.8 converges rather more rapidly since the coefficient of ϵ_n^2 is less. On the other hand, 9.3.9 requires only one division (to form $1/2b$) and this may be an advantage if it is to be used on an automatic digital calculator which, as is often the case, has no automatic divider.

Simple rearrangement of the original equation is not the only method for obtaining second order processes. Thus we may multiply the original equation $f(x) = 0$ by an arbitrary function $g(x)$ and apply the Newton-Raphson process to the combined system. If we take $g(x) = x$ and apply this to equation 9.3.7 there results:

$$x^3 - bx = 0$$

for which the iterative solution formula is:

$$x_{n+1} = x_n - \frac{x_n^3 - bx_n}{3x_n^2 - b} = 2x_n^3/(3x_n^2 - b). \quad \dots (9.3.10)$$

Now although both equations 9.3.8 and 9.3.10 are second order processes it is easy to show that, if:

$$x_n = \sqrt{b} + \epsilon$$

equation 9.3.8 leads to:

$$x_{n+1} = \sqrt{b} + \frac{1}{2}\epsilon^2/\sqrt{b} - \frac{1}{2}\epsilon^3/b + \dots$$

and equation 9.3.10 to:

$$x_{n+1} = \sqrt{b} + \frac{3}{2}\epsilon^2/\sqrt{b} - \frac{7}{2}\epsilon^3/b + \dots$$

Whence, by taking the iteration defined by:

$$\frac{1}{2}(9.3.8) - \frac{1}{2}(9.3.10)$$

the term in ϵ^2/\sqrt{b} is eliminated and we obtain the third order process:⁽³⁾

$$x_{n+1} = (5x_n^4 + 6bx_n^3 - 3b^2)/4x_n(3x_n^2 - b) + \epsilon^3/b + \dots \quad (9.3.11)$$

$$\text{Similarly from: } \frac{3}{4}(9.3.8) + \frac{1}{4}(9.3.9)$$

we obtain the third order process:

$$x_{n+1} = (3b^2 + 6bx_n^2 - x_n^4)/8bx_n - \frac{3}{4}\epsilon^3/b + \dots \quad (9.3.12)$$

and 9.3.11 and 9.3.12 could be combined to give a fourth order process. It is quite apparent from the complexity of 9.3.11 and 9.3.12 that they are of little utility in practical computation.

We may notice, before leaving the subject of second order processes for simple functions, the formula:

$$x_{n+1} = x_n(2 - bx_n) \quad \dots (9.3.13)$$

which converges to $1/b$ and is derived from the equation:

$$1/x - b = 0;$$

and also the iteration:

$$x_{n+1} = x_n[(p+1) - bx_n^p]/p \quad \dots (9.3.14)$$

which converges to $1/b^{1/p}$ and is derived from:

$$1/x^p - b = 0.$$

9.4 COMPLEX ROOTS

The Newton-Raphson process, described above, is usually considered to be applicable only to the determination of real roots, but actually it is quite feasible to use the Newton-Raphson process for evaluating complex roots, and since this is applicable to quite general equations and is second order, it will now be described.

Assume that it is desired to find a solution $z = a + ib$ to the equation:

$$f(z) = 0 \quad \dots (9.4.1)$$

Let $z = z_0 + \eta$, where $\eta = \epsilon + i\delta$, and ϵ and δ are assumed to be of the first order of small quantities. Then Taylor's theorem gives:

$$f(z) = 0 = f(z_0 + \eta) = f(z_0) + \eta f'(z_0) + \frac{1}{2}\eta^2 f''(z_0) + \dots$$

whence, to the second order:

$$\eta = -f(z_0)/f'(z_0) \quad \dots (9.4.2)$$

which is identical with the ordinary Newton-Raphson formula for real roots. We now note that:

$$\epsilon = \text{real part of } [-f(z_0)/f'(z_0)]$$

$$\delta = \text{imaginary part of } [-f(z_0)/f'(z_0)]$$

and that, to start the iteration, it is necessary to take a complex value of (z_0) which in the absence of other information can be taken to be (i).

As an example of the convergence of this process the following figures, obtained in solving:

$$z^2 - 2z + 2 = 0$$

from the initial value $z_0 = i$, show it in a favourable light:

$$\begin{array}{ll} z_0 = i & \eta_0 = \frac{1}{4}(3 - i) \\ z_1 = \frac{3}{4}(1 + i) & \eta_1 = \frac{1}{40}(13 + 9i) \\ z_2 = \frac{1}{40}(43 + 39i) & \eta_2 = -0.07673 + 0.02230i \\ \quad (= 1.075 + 0.975i) & \\ z_3 = 0.99827 + 0.99730i & \end{array}$$

The true solution is evidently $(1 + i)$ and the errors: $|1 + i - z_r|$ are 1, .35, .08 and .0032 respectively.

On the other hand, if an attempt is made to find directly a complex root of the equation:

$$z^3 - 3z^2 + 4z - 2 = 0$$

discussed in section 9.2, starting from the same initial approximation (i) the successive convergents are:

$$\begin{array}{ll} z_0 = i & \eta_0 = .46 - .24i \\ z_1 = .46 + .76i & \eta_1 = .43 - .12i \\ z_2 = .85 + .64i & \eta_2 = .825 + .28i \\ z_3 = 1.54 + .78i & \end{array}$$

with errors 1, .59, .39 and .58, so that z_3 is a less good approximation than z_2 . The reason for this behaviour lies in the fact that the derivative of

$$z^3 - 3z^2 + 4z - 2 = 0$$

$$\text{that is: } 3z^2 - 6z + 4$$

has a pair of complex roots $1 \pm .577i$ and that the value

$$z_2 = .85 + .64i$$

is fairly close to one of these.

The methods for determining complex roots of polynomial equations, which are usually recommended, are based upon the determination of the real quadratic factors of the polynomial. This avoids the need for complex arithmetic because the complex roots occur in conjugate pairs and, hence, imply a quadratic product with real coefficients.

Three methods for determining such real quadratic factors are in vogue: the method of iteration, usually known as Lin's method, the Newton-Raphson method, and Bairstow's process. Lin's method is first order and is based upon the following consideration. Suppose that our polynomial

$$f_n(z) = z^n + a_1 z^{n-1} + \dots + a_n = 0 \quad \dots (9.4.3)$$

is assumed to have a pair of complex, conjugate, roots α and $\bar{\alpha}$ which satisfy the quadratic equation:

$$q(z) = z^2 + b_1 z + b_2 = 0 \quad \dots (9.4.4)$$

Then

$$b_1 = -\alpha - \bar{\alpha} \quad \dots (9.4.5)$$

and

$$b_2 = \alpha \bar{\alpha} \quad \dots (9.4.6)$$

so that b_1 and b_2 are real.

If the quadratic 9.4.4 is an exact factor of 9.4.3 the result of dividing $f_n(z)$ by $q(z)$ will be:

$$f_n(z) = f_{n-2}(z) \cdot q(z)$$

where $f_{n-2}(z)$ is some polynomial, of degree $(n-2)$

$$f_{n-2}(z) = z^{n-2} + c_3 z^{n-3} + \dots + c_n \quad \dots (9.4.7)$$

In general $q(z)$ will be unknown; we therefore assume some trial quadratic $q_0(z) = z^2 + b_{10}z + b_{20}$ and attempt to divide $f_n(z)$ by it. The result can generally be written

$$f_n(z) = g_{n-2}(z)q_0(z) + p(z^2 + b_{11}z + b_{21})$$

where $g_{n-2}(z)$ is a polynomial of degree $(n-2)$ in z , of the type

$$g_{n-2}(z) = z^{n-2} + d_3 z^{n-3} + \dots + d_{n-2} z^3$$

In Lin's method the remainder polynomial $q_1(z) = z^2 + b_{11}z + b_{21}$ is used as a trial divisor and the process is repeated to obtain a sequence of quadratic divisors $q_i(z) = z^2 + b_{1i}z + b_{2i}$ until, if it converges, $q_{i+1} \rightarrow q_i$ at which point $q_i(z)$ is clearly a quadratic divisor of $f_n(z)$.

When a fairly good guess for the initial trial is available, Lin's method generally converges quite well, but convergence can be slow or non-existent and, for this reason, one of the methods now to be described is usually preferred.

The Newton-Raphson method for quadratic factors is second order, to derive it we assume that we are at the i th stage of an iterative process and that we have an approximate quadratic factor:

$$q_i(z) = z^2 + b_{1i}z + b_{2i} \quad \dots (9.4.8)$$

corresponding to 9.4.5 and 9.4.6 we have

$$b_{1i} = -\alpha_i - \bar{\alpha}_i \quad \dots (9.4.9)$$

$$b_{2i} = \alpha_i \bar{\alpha}_i \quad \dots (9.4.10)$$

Suppose that the true roots of 9.4.3 are $(\alpha_i + \delta\alpha_i)$, $(\bar{\alpha}_i + \delta\bar{\alpha}_i)$. The Newton-Raphson process, applied to 9.4.3, gives:

$$\delta\alpha_i = -f_n(\alpha_i)/f'_n(\alpha_i) \quad \dots (9.4.11)$$

$$\delta\bar{\alpha}_i = -f_n(\bar{\alpha}_i)/f'_n(\bar{\alpha}_i) \quad \dots (9.4.12)$$

whence, from 9.4.9 and 9.4.10, the changes required in b_{1i} and b_{2i} are given by:

$$\delta b_{1i} = -\delta\alpha_i - \delta\bar{\alpha}_i \quad \dots (9.4.13)$$

$$\delta b_{2i} = \bar{\alpha}_i \delta\alpha_i + \alpha_i \delta\bar{\alpha}_i \quad \dots (9.4.14)$$

To determine the various quantities involved we proceed as follows. Divide $f_n(z)$ by $q_i(z)$ to obtain

$$f_n(z) = q_i(z)f_{n-2}(z) + r_1 z + s_1 \quad \dots (9.4.15)$$

from which, since α_i and $\bar{\alpha}_i$ are the roots of $q_i(z) = 0$,

$$f_n(\alpha_i) = r_1 \alpha_i + s_1 \quad \dots (9.4.16)$$

$$f_n(\bar{\alpha}_i) = r_1 \bar{\alpha}_i + s_1 \quad \dots (9.4.17)$$

Next divide $f'_n(z)$ by $q_i(z)$, obtaining:

$$f'_n(z) = q_i(z)g_{n-2}(z) + r_2 z + s_2$$

whence, as before,

$$f'_n(\alpha_i) = r_2\alpha_i + s_2 \quad \dots (9.4.18)$$

$$f'_n(\bar{\alpha}_i) = r_2\bar{\alpha}_i + s_2 \quad \dots (9.4.19)$$

Substituting for 9.4.18–9.4.19 in 9.4.11 and 9.4.12 we get:

$$\delta\alpha_i = -\frac{r_1\alpha_i + s_1}{r_2\alpha_i + s_2}$$

$$\delta\bar{\alpha}_i = -\frac{r_1\bar{\alpha}_i + s_1}{r_2\bar{\alpha}_i + s_2}$$

and, substituting these values in 9.4.13 and 9.4.14, and using 9.4.9 and 9.4.10 to eliminate α_i and $\bar{\alpha}_i$, there results:

$$\delta b_{1i} = \frac{2r_1r_2b_{2i} - (r_1s_2 + r_2s_1)b_{1i} + 2s_1s_2}{r_2^2b_{2i} - r_2s_2b_{1i} + s_2^2} \quad \dots (9.4.20)$$

and

$$\delta b_{2i} = \frac{r_1r_2b_{1i}b_{2i} - 2r_1s_2b_{2i} + s_1s_2b_{1i} - r_2s_1(b_{1i}^2 - 2b_{2i})}{r_2^2b_{2i} - r_2s_2b_{1i} + s_2^2} \quad \dots (9.4.21)$$

We now define:

$$q_{i+1}(z) = z^2 + (b_{1i} + \delta b_{1i})z + (b_{2i} + \delta b_{2i}) = z^2 + b_{1i+1}z + b_{2i+1}$$

and repeat the process on $q_{i+1}(z)$.

The various quantities r_1, s_1, r_2, s_2 are formed by the usual process of synthetic division.

Bairstow's process is as follows. Assume, as before, that an i th stage approximation $q_i(z)$, defined by 9.4.8, is available. We again form the expression 9.4.15 obtaining $f_{n-2}(z)$, r_1 and s_1 . $f_{n-2}(z)$ is now divided by the trial quadratic giving:

$$f_{n-2}(z) = Q_i(z)f_{n-4}(z) + r_2z + s_2$$

The corrections δb_{1i} and δb_{2i} can then be shown⁽⁴⁾ to be:

$$\delta b_{1i} = \frac{r_1s_2 - r_2s_1}{r_2^2b_{2i} - r_2s_2b_{1i} + s_2^2}$$

$$\delta b_{2i} = \frac{r_1r_2b_{2i} - r_2s_1b_{1i} + s_1s_2}{r_2^2b_{2i} - r_2s_2b_{1i} + s_2^2}$$

A comparison of these equations with 9.4.20 and 9.4.21 shows that, when $f_n(z)$ is a quadratic, Bairstow's process gives an exact result whereas the Newton-Raphson method does not.

9.5 EQUAL ROOTS

When an equation is such that it is satisfied for values of the variable which tend to equality, with some parameter in the equation, the Newton-Raphson process will, in general, break down. Geometrically the situation is simply that the curve:

$$y = f(x)$$

is nearly parallel to the x axis in the region of the roots (or, in the limit touches the x axis), and this should have become evident in the preliminary survey.

When *exactly* equal roots occur they can be found by the simple process of extracting the highest common factor of $f(x)$ and $f'(x)$. When exact equality is not present, however, the first step should be to find the root of $f'(x) = 0$, say (x_m) . Then assume that the roots of $f(x) = 0$ are $(x_m \pm \epsilon)$ so that:

$$f(x_m \pm \epsilon) = 0 = f(x_m) \pm \epsilon f'(x_m) + \frac{1}{2}\epsilon^2 f''(x_m) \pm \dots$$

or, since by definition $f'(x_m) = 0$,

$$\epsilon = \pm \sqrt{\frac{-2f(x_m)}{f''(x_m)}} \quad \dots (9.5.1)$$

From this point it should be possible to revert to the normal Newton-Raphson process with some hope of convergence.

Alternatively, when a pair of equal roots is known to exist, an approximating value, x_0 , may be improved as follows. Assume that the correct solution is $(x_0 + \epsilon)$, then Taylor's theorem gives:

$$f(x_0 + \epsilon) = f(x_0) + \epsilon f'(x_0) + \frac{1}{2}\epsilon^2 f''(x_0) + \dots$$

Since a multiple root occurs at $(x_0 + \epsilon)$ this must correspond to the *minimum* of $f(x_0 + \epsilon)$ with respect to ϵ ,

$$\frac{df}{d\epsilon} = 0 \approx f'(x_0) + \epsilon f''(x_0)$$

whence: $\epsilon = -f'(x_0)/f''(x_0)$ (9.5.2)

9.6 ERRORS AND INACCURACIES

The solution of polynomial equations of high degree is, nowadays, almost always conducted on an automatic digital computer. The precision with which the calculations are to be conducted needs

careful examination and the following treatment, due to J. H. Wilkinson, is instructive.

Suppose that the polynomial equation is:

$$f(z) \equiv z^n + a_1 z^{n-1} + \dots + a_n = 0$$

and that we wish to examine the sensitivity of any root, say a , with respect to variations in one of the coefficients a_r . Let a_r become $(a_r + \delta a_r)$ and assume that the new solution is $(a + \delta a)$. We then have

$$f(a + \delta a) + \delta a_r (a + \delta a)^{n-r} = 0$$

Or, neglecting terms in δ^2 ,

$$f(a) + \delta a f'(a) + a^{n-r} \delta a_r = 0$$

But $f(a) = 0$ whence:

$$\frac{\delta a}{a} = -a^{n-r-1} \frac{\delta a_r}{f'(a)} = -a^{n-r-1} \frac{a_r}{f'(a)} \cdot \frac{\delta a_r}{a_r}$$

Wilkinson points out that if $n = 20$, and the quite well behaved polynomial $f(z) = (z + 1)(z + 2) \dots (z + 20)$ is considered, then for $a = -16$ the root is particularly sensitive to variations in the term a_1 . Thus

$$\frac{\delta a}{a} = -\frac{16^{18} \times 210}{4! \times 15!} \cdot \frac{\delta a_1}{a_1} \simeq 3.2 \times 10^{10} \frac{\delta a_1}{a_1}$$

so that calculations must be carried out to better than ten place accuracy.

The analysis points up the fact that, for equal roots (where $f'(a) = 0$), the solution to any equation is extremely sensitive to variations in any coefficient.

9.7 SIMULTANEOUS NON-LINEAR EQUATIONS

We return now to the solution of the general set of non-linear simultaneous equations envisaged in section 9.1

$$f_i(x_1, x_2, \dots, x_n) = 0 \quad (i = 1, \dots, m) \quad \dots (9.7.1)$$

that is, a set of m equations in n unknowns. When $m = n$ we have the classical case in which a finite multiplicity of solutions may exist. The more general situations, in which $m < n$ or $m > n$ are also encountered in physical problems and, in these circumstances there may be no *unique* solution in the former event, and no solution in the latter.

To unify the method of approach we shall define a solution of equations 9.7.1 to be any set of parameters (x_i) which makes the positive definite form:

$$\Phi = \sum_{i=1}^n G_i(f_i) \quad \dots (9.7.2)$$

a minimum, $G_i(z)$ being a function so chosen as to be always positive for all values of z , real or complex.

For convenience and simplicity it is usual to take either:

$$G_i(z) = z \bar{z}$$

or

$$G_i(z) = |z|$$

where \bar{z} is understood to mean the complex conjugate of z .

The general approach to the solution of a set of equations of type 9.7.1 should again be in two parts, first a survey to detect the rough location of the roots, and second a refinement process to evaluate the roots more exactly. Unfortunately the survey process is almost impossible when more than two variables are involved and in this event a refinement by successive approximation may have to be attempted.

When only a small number of equations and variables are involved (up to three) it may be feasible to use the Newton-Raphson process to find a solution to the required accuracy. Thus suppose that the required solution is at $(x_1 + \epsilon_1, x_2 + \epsilon_2, \dots, x_n + \epsilon_n)$, then the n variable form of Taylor's theorem gives:

$$\begin{aligned} f_i(x_1 + \epsilon_1, x_2 + \epsilon_2, \dots, x_n + \epsilon_n) \\ = f_i(x_1, x_2, \dots, x_n) + \sum_{r=1}^n \epsilon_r \frac{\partial f_i}{\partial x_r} + O(\epsilon^2) \quad (i = 1, \dots, n) \end{aligned}$$

whence, for a second order process, we take $(\epsilon_1, \dots, \epsilon_n)$ to be the solution of the set of *linear* simultaneous equations:

$$\epsilon_1 \frac{\partial f_i}{\partial x_1} + \epsilon_2 \frac{\partial f_i}{\partial x_2} + \dots + \epsilon_n \frac{\partial f_i}{\partial x_n} = -f_i \quad (i = 1, \dots, n) \quad \dots (9.7.3)$$

which correspond exactly to the one-dimensional Newton-Raphson process of equation 9.3.6.

When large numbers of variables and equations are involved this process will, in general, be too complicated to apply, and under these circumstances the analogue of one of the successive approximation methods of Chapter 7, sections 7.7. and 7.9 is to be preferred.

If the function $\Phi(x_1, x_2, \dots, x_n)$, defined in equation 9.7.2, is taken as the basis of a minimization process, the analysis can be conducted in a manner which closely parallels that of sections 7.7 and 7.9 with the exception that the n dimensional surfaces:

$$\Phi = \text{const.}$$

are no longer, in general, hyper-ellipsoids.

To take a two-dimensional example, the contour map of Figure 9.7.1 might be encountered. In this two maxima are shown at M_1

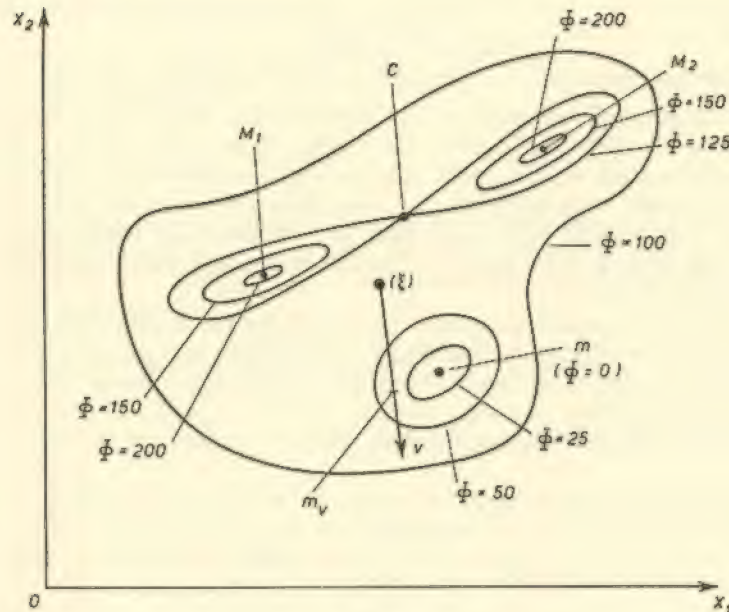


Figure 9.7.1

and M_2 , a minimum at m and a 'col', or pass between mountains, at C .

We seek to progress from an initial point of vector $\xi (= x_1 x_2 \dots x_n)$ to a vector $(\xi + a\mathbf{v})$ which makes $\Phi(\xi + a\mathbf{v})$ less than $\Phi(\xi)$. The Taylor expansion gives, at once,

$$\Phi(\xi + a\mathbf{v}) = \Phi(\xi) + \sum_{r=1}^n a v_r \frac{\partial \Phi}{\partial x_r} + \sum_{r,s=1}^n \frac{1}{2} a^2 v_r v_s \frac{\partial^2 \Phi}{\partial x_r \partial x_s} + O(a^3). \quad \dots (9.7.4)$$

Here \mathbf{v} is any vector and v_r is its component in direction x_r . Assume now that it is adequate to represent the variation of Φ along the vector \mathbf{v} by means of the terms of degree less than three in a in equation 9.7.4. We then obtain for the least value of Φ the multiplier a by means of:

$$\frac{d\Phi}{da} = 0.$$

That is, from equation 9.7.4,

$$a = - \frac{\sum_{r=1}^n v_r \frac{\partial \Phi}{\partial x_r}}{\sum_{r,s=1}^n v_r v_s \frac{\partial^2 \Phi}{\partial x_r \partial x_s}}. \quad \dots (9.7.5)$$

Just as in Chapter 7, there are several rational choices for the vector \mathbf{v} . If it is taken parallel to one of the axes of co-ordinates, x_k say, equation 9.7.5 becomes:

$$a = - \frac{\partial \Phi / \partial x_k}{\partial^2 \Phi / \partial x_k^2} \quad \dots (9.7.6)$$

and this corresponds to the relaxation technique of the linear case.

For a steepest descent, \mathbf{v} is defined by:

$$\mathbf{v} = \text{grad } \Phi = \sum_{r=1}^n \mathbf{e}_r \frac{\partial \Phi}{\partial x_r} \quad \dots (9.7.7)$$

where \mathbf{e}_r represents the unit vector parallel to the axis of x_r . Substituting in 9.7.5 we obtain:

$$a = - \frac{\sum_{r=1}^n \left(\frac{\partial \Phi}{\partial x_r} \right)^2}{\sum_{r,s=1}^n \frac{\partial \Phi}{\partial x_r} \cdot \frac{\partial \Phi}{\partial x_s} \cdot \frac{\partial^2 \Phi}{\partial x_r \partial x_s}} \quad \dots (9.7.8)$$

so that the actual change in any co-ordinate x_k is:

$$a \frac{\partial \Phi}{\partial x_k} = - \frac{\partial \Phi}{\partial x_k} \cdot \frac{\sum_{r=1}^n \left(\frac{\partial \Phi}{\partial x_r} \right)^2}{\sum_{r,s=1}^n \frac{\partial \Phi}{\partial x_r} \cdot \frac{\partial \Phi}{\partial x_s} \cdot \frac{\partial^2 \Phi}{\partial x_r \partial x_s}}. \quad \dots (9.7.9)$$

This formula is really too complex for use in actual computation and although it has been applied,⁽⁵⁾ is not so simple as a variant which will now be described.

Let the value of Φ along \mathbf{v} be as shown in Figure 9.7.2.

We first take the intersection of the tangent to the Φ , a curve with the line $\Phi = 0$, T say. OT is simply determined by means of equation 9.7.4 by neglecting terms in a^2 and so on. Thus:

$$a_T = -\Phi(\xi) / \sum_{r=1}^n v_r \frac{\partial \Phi}{\partial x_r} \quad \dots (9.7.10)$$

Next we calculate the values of Φ at the points $(\xi + a_T v)$ and $(\xi + \frac{1}{2}a_T v)$ and call these Φ_1 and $\Phi_{\frac{1}{2}}$ respectively.

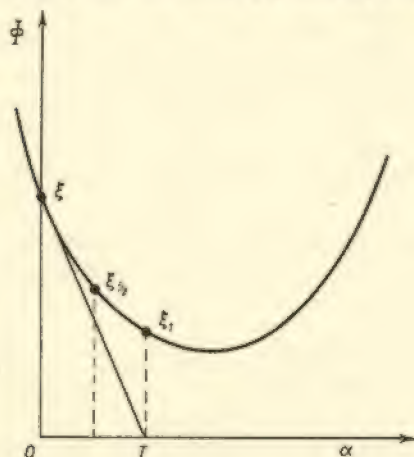


Figure 9.7.2

The minimum is then calculated by the simple process of forming the quadratic approximation to Φ which passes through Φ , $\Phi_{\frac{1}{2}}$ and Φ_1 , and then differentiating it with respect to α . Thus:

$$\Phi_a \approx \Phi + \frac{2a}{a_T} \Delta\Phi + \frac{a}{a_T} \left(\frac{2a}{a_T} - 1 \right) \Delta^2\Phi.$$

whence, since $\frac{d\Phi_a}{da} = 0$,

$$\frac{a}{a_T} = (\Delta^2\Phi - 2\Delta\Phi) / 4\Delta^2\Phi. \quad \dots (9.7.11)$$

Thus the individual component changes, av_r , are for the steepest descent, simply:

$$- \frac{(\Phi_1 - 4\Phi_{\frac{1}{2}} + 3\Phi)}{4(\Phi_1 - 2\Phi_{\frac{1}{2}} + \Phi)} \cdot \frac{\Phi \cdot \partial\Phi / \partial x_r}{\left[\sum_{r=1}^n \left(\frac{\partial\Phi}{\partial x_r} \right)^2 \right]}. \quad \dots (9.7.12)$$

It will be noticed that the use of this version of the steepest descent process requires the calculation of three values of Φ , and all $\frac{\partial\Phi}{\partial x_r}$,

$(n+3)$ values in all, as compared with about $\frac{n^2}{2}$ values of

$$\frac{\partial\Phi}{\partial x_r} \quad \text{and} \quad \frac{\partial^2\Phi}{\partial x_r \partial x_s}$$

for the analytic expression 9.7.9.

The method of minimizing Φ with respect to variations along a single vector v can be extended to cover the case of two vectors v and w . This requires the minimization of $\Phi(\xi + \alpha v + \beta w)$ with respect to α and β and follows the same lines as those used in deriving equation 9.7.5. Just as in Chapter 7, section 7.7 the vectors v and w are taken to be the current direction of steepest descent and the direction of the last descent; the formulae are, however, so complex as to be practically useless.

To conclude this section we may mention some of the difficulties which accompany the solution of non-linear algebraic equations by minimization methods. These stem from the fact that the surfaces $\Phi = \text{const.}$ are not simple hyper-ellipsoids having a single common centre corresponding to the minimum of Φ . As shown in Figure 9.7.1 maxima may occur, as at M_1 and M_2 , relative minima are possible, and 'cols' are of common occurrence. Each singularity of this type is a possible stationary value of Φ and care must be taken to ensure that any 'solution' obtained by descent methods is a true minimum and not one of these spurious addenda. This is particularly true in the case of redundant sets of equations where an actual zero of Φ is not to be hoped for.

It may be noted that, whilst maxima and 'cols' are easily detected by virtue of the fact that for the regions near to these the quadratic form:

$$\sum_{r,s} \frac{\partial^2\Phi}{\partial x_r \partial x_s} \cdot x_r x_s$$

is not positive definite, a relative minimum is quite indistinguishable by any purely mathematical test from the true minimum which defines the desired solution.

REFERENCES

- (1) OLVER, F. W. J., *Phil. Trans. Roy. Soc.*, A 244 (1952) 385
- (2) WHITTAKER, E. T. and ROBINSON, G., 'The Calculus of Observations,' 4th edn., p. 106 (1949)
- (3) HARTREE, D. R., *Proc. Camb. Phil. Soc.*, 45 (1948) 230
- (4) 'Modern Computing Methods,' p. 57, National Physical Laboratory, London. H.M. Stationery Office (1962)
- (5) BOOTH, A. D., *Quart. J. mech.*, II (1949) 460

APPROXIMATING FUNCTIONS

10.1 DEVELOPMENT IN SERIES

THE advent of automatic digital calculators has resulted in the frequent necessity for means of generating various common functions using only the operations of elementary arithmetic. Thus it is useless, from the computing machine viewpoint, to say 'calculate $\sin x$ ' since the machine will not normally understand a direct command of this type.

In practice a statement such as that made in the previous paragraph could be replaced by the equivalent:

$$\text{Form} \quad \left(x - \frac{x^3}{3!} + \frac{x^5}{5!} - \frac{x^7}{7!} \cdot \cdot \cdot \right)$$

taking sufficient terms so that the result is correct to the number of places required. In order to minimize the amount of work required to calculate the value of a given function it is thus appropriate to consider whether a better approximation than the Taylor series can be found for a given amount of work.

To make this more precise we may consider the following proposition: determine the $(N + 1)$ coefficients a_r so that the maximum difference between $f(x)$ and $\sum_{r=0}^N a_r x^r$ is as small as possible for a given range of values of x . Or again; determine the values of a_r so that the difference between $f(x)$ and $\sum_{r=0}^N a_r x^r$ has the least possible mean square value for a given range of x .

Which of these definitions will be best in any given context depends essentially upon the ultimate use to which the approximations to $f(x)$ are to be put. Clearly if only a single value is required the maximum divergence should be made as small as possible, on the other hand, if a number of such values are to be associated, either directly as:

$$\sum_{s=1}^M f(x_s) \quad \dots (10.1.1)$$

LEAST SQUARE APPROXIMATIONS

or with a weighting function $w(x_s)$

$$\sum_{s=1}^M w(x_s) f(x_s) \quad \dots (10.1.2)$$

then some approximation of least integral square deviation will be best.

In this chapter we shall consider several useful approximations of both categories.

10.2 LEAST SQUARE APPROXIMATIONS

Consider any set of functions $O_n(x)$ ($n = 0 \dots \infty$) ortho-normal in the interval (a, b) . We shall show that if an arbitrary function, $f(x)$, is expanded in a Fourier series of these functions, then the first $(n + 1)$ terms of that series constitute the best $(n + 1)$ term approximation to $f(x)$ in the least square sense.

For suppose that the expansion is:

$$f(x) \doteq \sum_{r=0}^n a_r O_r(x) \quad \dots (10.2.1)$$

then it is required that:

$$L = \int_a^b \left[f(x) - \sum_{r=0}^n a_r O_r(x) \right]^2 dx \quad \dots (10.2.2)$$

is a minimum. Now for a minimum:

$$\frac{\partial L}{\partial a_r} = 0 \quad (r = 0 \dots n) \quad \dots (10.2.3)$$

whence:

$$\int_a^b O_s(x) \left[f(x) - \sum_{r=0}^n a_r O_r(x) \right] dx = 0 \quad (s = 0 \dots n) \quad \dots (10.2.4)$$

but, since the functions $O_r(x)$ are ortho-normal in (a, b)

$$\int_a^b O_r(x) O_s(x) dx = \delta_{rs} \quad (\dots 10.2.5)$$

where δ_{rs} is the Kronecker delta function defined by:

$$\delta_{rs} = 0 (r \neq s) \\ \delta_{rs} = 1 (r = s)$$

APPROXIMATING FUNCTIONS

whence equation 10.2.4 may be written:

$$\int_a^b O_s(x)f(x)dx = \sum_{r=0}^n a_r \delta_{rs} = a_s \quad \dots (10.2.6)$$

but this is simply the definition of the coefficient of $O_s(x)$ in the Fourier expansion of $f(x)$.

The approximation can be generalized to include a weight function $w(x)$. Thus the functions $O_n(x)$ satisfy:

$$\int_a^b O_r(x)O_s(x)w(x)dx = \delta_{rs} \quad \dots (10.2.7)$$

and

$$L_w = \int_a^b w(x) \left[f(x) - \sum_{r=0}^n a_r O_r(x) \right]^2 dx \quad \dots (10.2.8)$$

is a minimum when

$$a_s = \int_a^b w(x) O_s(x) f(x) dx \quad \dots (10.2.9)$$

which is, again, the normal Fourier coefficient.

10.3 SOME USEFUL FUNCTIONS FOR LEAST SQUARE APPROXIMATION

From the classical viewpoint the most important set of orthogonal functions are $\sin(nx)$ and $\cos(nx)$, which lead to the well-known Fourier expansions in $(0, \pi)$ and $(0, 2\pi)$. These will not be further considered here save to remark that, under suitable circumstances, not only is the mean square difference between function and approximation minimized, but that the same is true of derivatives.

The most important polynomial expansions, for unit weight function, are these in terms of the Legendre polynomials $P_n(x)$. These are defined by:

$$P_n(x) = \frac{1}{2^n n!} \frac{d^n}{dx^n} (x^2 - 1)^n \quad \dots (10.3.1)$$

$$P_{n+1}(x) = [(2n+1)xP_n(x) - nP_{n-1}(x)]/(n+1) \quad \dots (10.3.2)$$

with orthogonality relations:

$$\int_{-1}^{+1} P_n(x)P_m(x)dx = 2\delta_{mn}/(2n+1). \quad \dots (10.3.3)$$

USEFUL FUNCTIONS FOR LEAST SQUARE APPROXIMATION

For intervals other than $(-1, +1)$ it is easy to transform the Legendre polynomials by means of:

$$x = \frac{2}{b-a}z - \frac{b+a}{b-a} \quad \dots (10.3.4)$$

which gives a set of polynomials $Q_n(z)$, say, orthogonal in (a, b) and with orthogonality relations:

$$\int_a^b Q_n(z)Q_m(z)dz = (b-a)\delta_{mn}/(2n+1). \quad \dots (10.3.5)$$

The first eleven Legendre polynomials for interval $(-1, +1)$ are given in Table 10.3.1.

Table 10.3.1. Legendre polynomials for interval $(-1, +1)$

n	$P_n(x)$	n	$P_n(x)$
0	1	6	$\frac{1}{8}(231x^6 - 315x^4 + 105x^2 - 5)$
1	x	7	$\frac{1}{8}(429x^7 - 693x^5 + 315x^3 - 35x)$
2	$\frac{1}{2}(3x^2 - 1)$	8	$\frac{1}{16}(6435x^8 - 12012x^6 + 6930x^4 - 1260x^2 + 35)$
3	$\frac{1}{2}(5x^3 - 3x)$	9	$\frac{1}{16}(12155x^9 - 25740x^7 + 18018x^5 - 4620x^3 + 315x)$
4	$\frac{1}{8}(35x^4 - 30x^2 + 3)$	10	$\frac{1}{32}(46189x^{10} - 109395x^8 + 90090x^6 - 30030x^4 + 3465x^2 - 63)$
5	$\frac{1}{8}(63x^5 - 70x^3 + 15x)$		

Useful expansions are:

$$x^{2n} = \frac{1}{(2n+1)} \left[1 \cdot P_0(x) + 5 \cdot \frac{2n}{2n+3} \cdot P_2(x) + 9 \cdot \frac{2n(2n-2)}{(2n+3)(2n+5)} \cdot P_4(x) + \dots \right]$$

$$x^{2n+1} = \frac{1}{(2n+3)} \left[3 \cdot P_1(x) + 7 \cdot \frac{2n}{2n+5} P_3(x) + 11 \cdot \frac{2n(2n-2)}{(2n+5)(2n+7)} P_5(x) + \dots \right]$$

$$(1-x^2)^{\frac{1}{2}} = \frac{\pi}{2} \left[\frac{1}{2} P_0(x) - 5 \cdot \frac{1}{4} \left(\frac{1}{2} \right)^2 P_2(x) - 9 \cdot \frac{3}{8} \left(\frac{1}{2.4} \right)^2 P_4(x) - 13 \cdot \frac{5}{8} \left(\frac{1.3}{2.4.6} \right)^2 P_6(x) - \dots \right]$$

$$(1-x^2)^{-\frac{1}{2}} = \frac{\pi}{2} \left[P_0(x) + 5 \cdot \left(\frac{1}{2}\right)^2 P_2(x) + 9 \cdot \left(\frac{1.3}{2.4}\right)^2 P_4(x) + 13 \cdot \left(\frac{1.3.5}{2.4.6}\right)^2 P_6(x) + \dots \right]$$

$$\sin^{-1}(x) = \frac{\pi}{8} \left[3 \cdot P_1(x) + 7 \cdot \left(\frac{1}{4}\right)^2 P_3(x) + 11 \cdot \left(\frac{1.3}{4.6}\right)^2 P_5(x) + \dots \right]$$

$$\sin\left(\frac{\pi}{2}x\right) = 3 \cdot J_{\frac{1}{2}}(\pi/2)P_1(x) - 7 \cdot J_{\frac{3}{2}}(\pi/2)P_3(x) + 11 \cdot J_{\frac{5}{2}}(\pi/2)P_5(x) - \dots$$

$$\cos\left(\frac{\pi}{2}x\right) = 1 \cdot J_{\frac{1}{2}}(\pi/2)P_0(x) - 5 \cdot J_{\frac{3}{2}}(\pi/2)P_2(x) + 9 \cdot J_{\frac{5}{2}}(\pi/2)P_4(x) - \dots$$

where $J_\nu(x)$ is the ν th order Bessel function.

Polynomials appropriate to $(0, \infty)$ and for weight function e^{-x} are those of Laguerre, they are given by:

$$L_n(x) = e^x \frac{d^n}{dx^n} (x^n e^{-x}) \quad \dots (10.3.6)$$

$$L_{n+1}(x) = (2n+1-x)L_n(x) - n^2 L_{n-1}(x) \quad \dots (10.3.7)$$

and satisfy:

$$\int_0^\infty L_n(x) L_m(x) e^{-x} dx = (n!)^2 \delta_{mn}. \quad \dots (10.3.8)$$

Finally, for interval $(-\infty, +\infty)$, there are Hermite polynomials, which are defined by:

$$H_n(x) = (-1)^n e^{x^2} \frac{d^n}{dx^n} (e^{-x^2}) \quad \dots (10.3.9)$$

$$H_{n+1}(x) = 2xH_n(x) - 2nH_{n-1}(x) \quad \dots (10.3.10)$$

and have orthogonality relations:

$$\int_{-\infty}^\infty H_n(x) H_m(x) e^{-x^2} dx = 2^n \cdot n! \sqrt{\pi} \delta_{mn}. \quad \dots (10.3.11)$$

The expressions for some of the above polynomials are given in Table 10.3.2.

Table 10.3.2. Laguerre and Hermite polynomials

n	$L_n(x)$	$H_n(x)$
0	1	1
1	$-x+1$	$2x$
2	x^2-4x+2	$4x^2-2$
3	$-x^3+9x^2-18x+6$	$8x^3-12x$
4	$x^4-16x^3+72x^2-96x+24$	$16x^4-48x^2+12$
5	$-x^5+25x^4-200x^3+600x^2-600x+120$	$32x^5-160x^3+120x$

10.4 LEAST ABSOLUTE DEVIATION APPROXIMATION

A remarkable set of orthogonal polynomials is that derived by Chebyshev. These are defined by:

$$T_0^* = 1, \quad T_n^*(x) = \frac{1}{2^{n-1}} \cos(n \cos^{-1} x) \quad (n \geq 1) \quad \dots (10.4.1)$$

$$T_{n+1}^*(x) = xT_n^*(x) - \frac{1}{2}T_{n-1}^*(x) \quad \dots (10.4.2)$$

and have orthogonality relations:

$$\int_{-1}^{+1} T_n^*(x) T_m^*(x) \frac{dx}{\sqrt{1-x^2}} = \frac{\pi}{2^{2n-1}} \delta_{mn}. \quad \dots (10.4.3)$$

Their importance lies in the fact that, of all polynomials of degree n , $T_n^*(x)$ has the least maximum value in $(-1, +1)$ when the coefficient of x^n is taken to be unity. For consider the function defined by equation 10.4.1, it is easily seen that the maxima and minima occur at:

$$x_k = \cos k\pi/n$$

where $k = 0, 1, \dots, n$, and that for these points:

$$T_n^*(x_k) = (-1)^{k/2^{n-1}}.$$

Let

$$M_n(x) = x^n + a_1 x^{n-1} + a_2 x^{n-2} + \dots + a_n$$

be a function, not everywhere zero, whose deviation from zero in $(-1, +1)$ is, if possible, less than that of $T_n^*(x)$. Then:

$$T_n^*(x_0) - M_n(x_0) > 0, \quad T_n^*(x_1) - M_n(x_1) < 0, \quad T_n^*(x_2) - M_n(x_2) > 0 \text{ etc.}$$

that is, the function $T_n^*(x) - M_n(x)$ has roots between (x_0, x_1) , (x_1, x_2) , ..., (x_{n-1}, x_n) . There are n of these, and since $T_n^* - M_n$ is of degree $(n-1)$ at most, we are led to a contradiction, it follows that the assumption that $M_n(x)$ is different from zero and $T_n^*(x)$ is false.

Modern usage⁽¹⁾ favours the slightly modified definition:

$$T_0 = 1 \quad T_n(x) = \cos(n \cos^{-1} x) \quad \dots (10.4.4)$$

$$T_{n+1} = 2xT_n(x) - T_{n-1}(x) \quad \dots (10.4.5)$$

$$\left. \begin{aligned} \int_{-1}^{+1} T_n(x) T_m(x) \frac{dx}{\sqrt{1-x^2}} &= \frac{\pi}{2} \delta_{mn} \quad (m \neq 0) \\ &= \pi \quad (m = n = 0) \end{aligned} \right\} \dots (10.4.6)$$

which has the advantage of avoiding fractions in the values of $T_n(x)$, the first eleven of which are given in *Table 10.4.1*.

Table 10.4.1. Chebyshev polynomials

n	$T_n(x)$	n	$T_n(x)$
0	1	6	$32x^6 - 48x^4 + 18x^2 - 1$
1	x	7	$64x^7 - 112x^5 + 56x^3 - 7x$
2	$2x^2 - 1$	8	$128x^8 - 256x^6 + 160x^4 - 32x^2 + 1$
3	$4x^3 - 3x$	9	$256x^9 - 576x^7 + 432x^5 - 120x^3 + 9x$
4	$8x^4 - 8x^2 + 1$	10	$512x^{10} - 1280x^8 + 1120x^6 - 400x^4 + 50x^2 - 1$
5	$16x^5 - 20x^3 + 5x$		

Useful expansions are:

$$x^n = \frac{1}{2^{n-1}} \sum_{k=0}^{[n]} \binom{n}{k} T_{n-2k}(x)$$

where $[n]$ represents the largest integer in $(\frac{1}{2}n)$ and the coefficient of $T_0(x)$ (if present) is halved,

$$\text{and:} \quad \sin\left(\frac{\pi}{2}x\right) = 2 \sum_{r=0}^{\infty} (-1)^r J_{2r+1}\left(\frac{\pi}{2}\right) \cdot T_{2r+1}(x)$$

$$\cos\left(\frac{\pi}{2}x\right) = J_0\left(\frac{\pi}{2}\right) + 2 \sum_{r=1}^{\infty} (-1)^r J_{2r}\left(\frac{\pi}{2}\right) \cdot T_{2r}(x)$$

where $J_r(x)$ is the r th order Bessel function.

These Bessel functions, and also those required in the Legendre Polynomial expansions, do not appear to be readily available; we give their values in *Table 10.4.2*.

Table 10.4.2. Bessel functions for Legendre and Chebyshev interpolation

n	$J_n(\pi/2)$	$J(n + \frac{1}{2})(\pi/2)$
0	0.47200 12157 7	0.63661 97723 7
1	0.56682 40889 1	0.40528 47345 7
2	0.24970 16291 4	0.13741 70540 3
3	0.06903 58882 9	0.03212 73337 1
4	0.01399 60398 1	0.00575 32170 8
5	0.00224 53571 2	0.00083 61720 0
6	0.00029 83476 0	0.00010 23428 0
7	0.00003 38506 4	0.00001 08228 5
8	0.00000 33522 0	0.00000 10077 8
9	0.00000 02945 7	0.00000 00838 4
10	0.00000 00232 7	0.00000 00063 0
11	0.00000 00016 7	0.00000 00004 3
12	0.00000 00001 1	0.00000 00000 3

10.5 GENERATING FUNCTIONS AND DIFFERENTIAL EQUATIONS

In deriving expansions using the various polynomials discussed above it is sometimes convenient to have them expressed in terms of a generating function. The most useful of these are:

Legendre Polynomials:

$$\frac{1}{\sqrt{1-2hx+h^2}} = \sum_{n=0}^{\infty} h^n P_n(x) \quad \dots (10.5.1)$$

Laguerre Polynomials:

$$\frac{e^{-xh/(1-h)}}{(1-h)} = \sum_{n=0}^{\infty} \frac{h^n}{n!} L_n(x) \quad \dots (10.5.2)$$

Hermite Polynomials:

$$e^{[x^2-(h-x)^2]} = \sum_{n=0}^{\infty} \frac{h^n}{n!} H_n(x) \quad \dots (10.5.3)$$

Chebyshev Polynomials:

$$\frac{1-xh}{1-2xh+h^2} = \sum_{n=0}^{\infty} h^n T_n(x) \quad \dots (10.5.4)$$

APPROXIMATING FUNCTIONS

The differential equations satisfied by P_n , L_n , H_n and T_n are ⁽²⁾, respectively:

$$(x^2 - 1)P_n''(x) + 2xP_n'(x) - n(n+1)P_n(x) = 0 \quad \dots (10.5.5)$$

$$xL_n''(x) + (1-x)L_n'(x) + nL_n(x) = 0 \quad \dots (10.5.6)$$

$$H_n''(x) - 2xH_n'(x) + 2nH_n(x) = 0 \quad \dots (10.5.7)$$

and

$$(1-x^2)T_n''(x) - xT_n'(x) + n^2T_n(x) = 0. \quad \dots (10.5.8)$$

10.6 A COMPARISON OF ACCURACY

We conclude this brief account of the use of approximating polynomials by giving a table, due to GOODWIN ⁽²⁾, showing the relative number of terms required for the ordinary Taylor series and of the Chebyshev approximating polynomial in certain cases.

Function	Range	Terms in Taylor expansion	Terms in Chebyshev expansion
e^x	$(0, +1)$	14	9
$\cos x$	$(-\pi/2, \pi/2)$	8	7
$\sin^{-1} x$	$(-1, +1)$	25	10
$\ln(1+x)$	$(0, +1)$	10^{10}	14

In each case the required accuracy is ten decimal places.

REFERENCES

- ⁽¹⁾ LANCZOS, C., 'Tables of Chebyshev Polynomials,' *Nat. Bur. Stand. Appl. Maths. Series, No. 9*, Washington (1952)
- ⁽²⁾ COURANT, R. and HILBERT, D., 'Methods of Mathematical Physics,' vol. 1, Chap. II, New York (1953)
- ⁽³⁾ GOODWIN, E. T., 'The Use of Mathematical Tables in a High-speed Digital Computer,' N.P.L. Symposium on automatic digital computation. Paper No. 21. Teddington (1953)

FOURIER SYNTHESIS AND ANALYSIS

11.1 FOURIER SYNTHESIS

ONE major source of numerical work in modern physics lies in the deduction of molecular structure from observations of the x-ray diffraction spectra which arise when a monochromatic beam of x-rays is incident, at the appropriate BRAGG ⁽¹⁾ angle, on a single crystal. This leads to a Fourier synthesis, a type of calculation which occurs also in tide prediction and in astrophysics.

For the purposes of the numerical analyst it is sufficient to note that Fourier synthesis calls for the evaluation of functions $\rho(x, y, z)$ defined by:

$$\rho(x, y, z) = \frac{1}{V} \sum_h \sum_k \sum_l |F(h, k, l)| \cos \left[2\pi \left(h \frac{x}{a} + k \frac{y}{b} + l \frac{z}{c} \right) - \alpha_{hkl} \right]. \quad \dots (11.1.1)$$

Typical cases have 60 values of x, y and z , and thus involve the evaluation of $(60)^3$ values of $\rho(x, y, z)$. The observed coefficients $|F(h, k, l)|$ are frequently up to 1000 in number and, in some recent work, have been as many as 10,000.

Under these circumstances it is not practicable, even with a high speed automatic digital calculator, to evaluate the $\rho(x, y, z)$ individually, and means have been devised for the simultaneous evaluation of whole groups of $\rho(x, y, z)$; usually for fixed values of two of the variables and a range of values of the third.

A very large number of different summation methods have been evolved, but we shall confine our attention to the method in most common use, which is due to BEEVERS and LIPSON. ^(2, 3)

The Beevers-Lipson method depends upon the breaking down of the three dimensional synthesis (11.1.1) into a series of one dimensional syntheses of the type:

$$\rho(x) = \sum_{h=-H_1}^{+H_1} \pm A_h \begin{Bmatrix} \sin \\ \cos \end{Bmatrix} \left(2\pi h \frac{x}{a} \right) \quad \dots (11.1.2)$$

which, as we shall presently show, can be done.

To effect the summation (11.1.2), for a range of values of (x) , Beevers and Lipson prepared a set of 'strips', one to each value of

A_h in the integer range 1-99 and for each value of h in the range 0-20. The sine and cosine functions are catered for separately. Each strip bears, upon its front surface, the information:

$$\begin{array}{ccc} \text{(Value of } A) & \text{(Sine or Cosine)} & \text{(Frequency } h) \\ (16 \text{ Values of } A \begin{cases} \sin \\ \cos \end{cases} \frac{2\pi}{60} n \cdot h) \end{array}$$

where n runs through the integers 0-15; the reverse side of the strip contains the same information, but for $-A$. A typical strip is shown in Figure 11.1.1.

Figure 11.1.1 Beavers-Lipson strip for $32 \sin 2\pi \cdot 3(x/a)$

32S3	0	10	19	26	30	32	30	26	19	10	0	10	19	26	30	32
------	---	----	----	----	----	----	----	----	----	----	---	----	----	----	----	----

The reason for the choice of interval $2\pi/60$ is a physical one which is related to the resolving power to be expected from the use of x-rays of a given wavelength.

The use of the strips will be clear from the following simple example, suppose that we require:

$$\mathcal{E} = 12 \cos 2\pi \cdot 1(x/a) + 32 \sin 2\pi \cdot 3(x/a) - 40 \cos 2\pi \cdot 4(x/a). \quad \dots(11.1.3)$$

The appropriate strips are selected and laid together in a matrix:

$x/a =$	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
12 C 1	12	12	12	11	11	10	10	9	8	7	6	5	4	3	1	0
32 S 3	0	10	19	26	30	32	30	26	19	10	0	10	19	26	30	32
40 C 4	40	37	27	12	4	20	32	39	39	32	20	4	12	27	37	40
Σ	28	15	4	25	45	62	72	74	66	49	26	1	27	50	68	72

If it is desired to extend the range beyond $(\frac{15}{60} \times 2\pi)$ it is merely necessary to reverse the cosine strips of *odd* frequency and the sine strips of *even* frequency, thus in the above case:

$x/a =$	30	29	28	27	26	25	24	23	22	21	20	19	18	17	16	15
12 C 1	12	12	12	11	11	10	10	9	8	7	6	5	4	3	1	0
32 S 3	0	10	19	26	30	32	30	26	19	10	0	10	19	26	30	32
40 C 4	40	37	27	12	4	20	32	39	39	32	20	4	12	27	37	40
Σ	62	39	20	3	23	42	52	56	50	35	14	11	35	56	68	72

The ranges (30-45) (45-60) follow in a similar manner from the symmetry properties of the trigonometric functions. In more complicated summations it is more economical to perform the sine-even, sine-odd, cosine-even, cosine-odd summations separately and then to combine the results as required to extend the range.

We now revert to the original summation (11.1.1). This may be written:

$$\begin{aligned} \rho(x, y, z) = \frac{1}{V} \sum_{h,k,l} \sum \sum |F| \cos a \cos 2\pi \left(h \frac{x}{a} + k \frac{y}{b} + l \frac{z}{c} \right) \\ + |F| \sin a \sin 2\pi \left(h \frac{x}{a} + k \frac{y}{b} + l \frac{z}{c} \right) \quad \dots(11.1.4) \end{aligned}$$

where, for simplicity, we have written $|F|$ for $|F(h, k, l)|$ and a for a_{hkl} .

In a like manner 11.1.4 may be broken down into:

$$\begin{aligned} \rho(x, y, z) = \frac{1}{V} \sum_{h,k} \sum \cos 2\pi \left(h \frac{x}{a} + k \frac{y}{b} \right) \sum_l \left[|F| \cos a \cos 2\pi l \frac{z}{c} + |F| \sin a \sin 2\pi l \frac{z}{c} \right] \\ + \sin 2\pi \left(h \frac{x}{a} + k \frac{y}{b} \right) \sum_l \left[|F| \sin a \cos 2\pi l \frac{z}{c} - |F| \cos a \sin 2\pi l \frac{z}{c} \right] \quad \dots(11.1.5) \end{aligned}$$

which is of the form:

$$\rho(x, y, z) = \frac{1}{V} \sum_{h,k} \sum C_z \cos 2\pi \left(h \frac{x}{a} + k \frac{y}{b} \right) + S_z \sin 2\pi \left(h \frac{x}{a} + k \frac{y}{b} \right) \quad \dots(11.1.6)$$

where:

$$C_z = \sum_l \left[|F| \cos a \cos 2\pi l \frac{z}{c} + |F| \sin a \sin 2\pi l \frac{z}{c} \right] \quad \dots(11.1.7)$$

$$S_z = \sum_l \left[|F| \sin a \cos 2\pi l \frac{z}{c} - |F| \cos a \sin 2\pi l \frac{z}{c} \right] \quad \dots(11.1.8)$$

Now the summations for C_z and S_z can be performed directly by means of the Beavers-Lipson strips and the resulting sets of coefficients tabulated. The process is then repeated to give:

$\rho(x, y, z) =$

$$\frac{1}{V} \sum_h \cos 2\pi \left(h \frac{x}{a} \right) \sum_k \left[h k C_z \cos 2\pi \left(k \frac{y}{b} \right) + h k S_z \sin 2\pi \left(k \frac{y}{b} \right) \right] \\ + \sin 2\pi \left(h \frac{x}{a} \right) \sum_k \left[h k S_z \cos 2\pi \left(k \frac{y}{b} \right) - h k C_z \sin 2\pi \left(k \frac{y}{b} \right) \right]$$

or

$$\rho(x, y, z) = \frac{1}{V} \sum_h C_{yz} \cos 2\pi \left(h \frac{x}{a} \right) + h S_{yz} \sin 2\pi \left(h \frac{x}{a} \right) \dots (11.1.9)$$

where:

$$h C_{yz} = \sum_k \left[h k C_z \cos 2\pi \left(k \frac{y}{b} \right) + h k S_z \sin 2\pi \left(k \frac{y}{b} \right) \right] \dots (11.1.10)$$

$$h S_{yz} = \sum_k \left[h k S_z \cos 2\pi \left(k \frac{y}{b} \right) - h k C_z \sin 2\pi \left(k \frac{y}{b} \right) \right] \dots (11.1.11)$$

Thus, since $h C_{yz}$, $h S_{yz}$ can be evaluated by the Beevers-Lipson technique, and the final summation 11.1.9 is also in the correct form, the problem is solved.

Various points of technique arise in the practical summation of multiple series of the type 11.1.4, notably that of 'multiplicity correction'. This is necessary because the method of reduction used in 11.1.5-11.1.11 has the effect of including terms of the types $(h, k, 0)$ $(h, 0, 0)$ $(0, 0, 0)$ more than once. It is sufficient to state here that the corrections:

$$\begin{aligned} F(h, k, l) &\div 1 \\ F(h, k, 0) &\div 2 \\ F(h, 0, 0) &\div 4 \\ F(0, 0, 0) &\div 8 \end{aligned}$$

are appropriate, for further details the reader is referred to the specialist monographs.^(4, 5)

11.2 THE LOCATION OF MAXIMA

One important requirement in dealing with series of the type envisaged in section 11.1 is that of locating accurately the maxima of $\rho(x, y, z)$. This can be conveniently achieved by the methods of differential synthesis^(6, 7) which make unnecessary a complete

evaluation of $\rho(x, y, z)$ for all values of (x, y, z) . The method is essentially that of Newton-Raphson (see Chapter 9, section 9.3) applied to the derivatives of $\rho(x, y, z)$. Thus, at a maximum we have:

$$\frac{\partial \rho}{\partial x} = \frac{\partial \rho}{\partial y} = \frac{\partial \rho}{\partial z} = 0.$$

We denote the synthesis 11.1.1 by the shorthand notation:

$$\rho = \frac{1}{V} \sum_3 |F| \cos (\theta - \alpha) \dots (11.2.1)$$

whence

$$\left. \begin{aligned} \frac{\partial \rho}{\partial x} &= -\frac{2\pi}{aV} \sum_3 |F| \sin (\theta - \alpha) = 0 \\ \frac{\partial \rho}{\partial y} &= -\frac{2\pi}{bV} \sum_3 k |F| \sin (\theta - \alpha) = 0 \\ \frac{\partial \rho}{\partial z} &= -\frac{2\pi}{cV} \sum_3 l |F| \sin (\theta - \alpha) = 0 \end{aligned} \right\} \dots (11.2.2)$$

Assume that (x, y, z) is an approximation to the maximum and that the true value is $(x + \epsilon_x, y + \epsilon_y, z + \epsilon_z)$, then, to the first order:

$$\left. \begin{aligned} -\frac{2\pi}{aV} \sum_3 \{h |F| \sin (\theta - \alpha) + h \epsilon_\theta \cos (\theta - \alpha)\} &= 0 \\ -\frac{2\pi}{bV} \sum_3 \{k |F| \sin (\theta - \alpha) + k \epsilon_\theta \cos (\theta - \alpha)\} &= 0 \\ -\frac{2\pi}{cV} \sum_3 \{l |F| \sin (\theta - \alpha) + l \epsilon_\theta \cos (\theta - \alpha)\} &= 0 \end{aligned} \right\} \dots (11.2.3)$$

where

$$\epsilon_\theta = 2\pi \left(h \frac{\epsilon_x}{a} + k \frac{\epsilon_y}{b} + l \frac{\epsilon_z}{c} \right).$$

From these equations we obtain:

$$\left. \begin{aligned} A_{hh} \epsilon_x + A_{hk} \epsilon_y + A_{hl} \epsilon_z + A_h &= 0 \\ A_{kh} \epsilon_x + A_{kk} \epsilon_y + A_{kl} \epsilon_z + A_k &= 0 \\ A_{lh} \epsilon_x + A_{lk} \epsilon_y + A_{ll} \epsilon_z + A_l &= 0 \end{aligned} \right\} \dots (11.2.4)$$

where:

$$\left. \begin{aligned} A_h &= -\frac{2\pi}{aV} \sum_3 h |F| \sin(\theta - \alpha) \\ A_k &= -\frac{2\pi}{bV} \sum_3 k |F| \sin(\theta - \alpha) \\ A_l &= -\frac{2\pi}{cV} \sum_3 l |F| \sin(\theta - \alpha) \end{aligned} \right\} \dots (11.2.5)$$

$$\left. \begin{aligned} A_{hk} &= A_{kh} = -\frac{4\pi^2}{abV} \sum_3 hk |F| \cos(\theta - \alpha) \\ A_{kl} &= A_{lk} = -\frac{4\pi^2}{bcV} \sum_3 kl |F| \cos(\theta - \alpha) \\ A_{lh} &= A_{hl} = -\frac{4\pi^2}{caV} \sum_3 lh |F| \cos(\theta - \alpha) \end{aligned} \right\} \dots (11.2.6)$$

$$\left. \begin{aligned} A_{hh} &= -\frac{4\pi^2}{a^2V} \sum_3 h^2 |F| \cos(\theta - \alpha) \\ A_{kk} &= -\frac{4\pi^2}{b^2V} \sum_3 k^2 |F| \cos(\theta - \alpha) \\ A_{ll} &= -\frac{4\pi^2}{c^2V} \sum_3 l^2 |F| \cos(\theta - \alpha) \end{aligned} \right\} \dots (11.2.7)$$

The process as described above is second order, in practice it is usual to compute the coefficients A_{hh} , A_{kk} , A_{ll} , A_{hk} , A_{kl} , A_{lh} only at the start and at the finish of the refinement which makes the process formally only first order; the convergence is, however, still extremely rapid.

11.3 RADIAL AND OTHER SYNTHESSES

For practical requirements it is sometimes necessary to obtain from the three dimensional x-ray data a radial distribution function. This is defined as the average of $\rho(x, y, z)$ over shells of constant radius from the origin. The series involved is usually:

$$P(u, v, w) = \frac{1}{V} \sum_{h,k,l} \sum \sum |F(h, k, l)|^2 \cos 2\pi \left(h \frac{u}{a} + k \frac{v}{b} + l \frac{w}{c} \right) \dots (11.3.1)$$

the so-called PATTERSON function⁽⁸⁾, or to mathematicians the 'convolution' or 'faltung',⁽⁹⁾ and the radial distribution function $p_3(r)$ is defined to be:

$$4\pi r^2 p_3(r) dr = \int_S P(u, v, w) dS \cdot dr \dots (11.3.2)$$

where dS is an element of the spherical surface, of radius r , whose centre is $(0, 0, 0)$.

Now the expression 11.3.2 may be written, by virtue of equation 11.3.1,

$$4\pi r^2 p_3(r) dr = \frac{1}{V} \sum_{h,k,l} \sum \sum |F(h, k, l)|^2 \int_S \cos 2\pi \left(h \frac{u}{a} + k \frac{v}{b} + l \frac{w}{c} \right) dS \cdot dr$$

and the integrals can be evaluated by transforming to axes normal and perpendicular to the plane:

$$h \frac{u}{a} + k \frac{v}{b} + l \frac{w}{c} = 0.$$

Thus:

$$\begin{aligned} \int_S \cos 2\pi \left(h \frac{u}{a} + k \frac{v}{b} + l \frac{w}{c} \right) dS &= \int_{-r}^{+r} \cos 2\pi \left(\frac{n}{d(h, k, l)} \right) \cdot 2\pi r \cdot dn \\ &= 4\pi r^2 \frac{\sin 2\pi r/d(h, k, l)}{2\pi r/d(h, k, l)} \dots (11.3.3) \end{aligned}$$

where dn is an element in the direction of the normal to

$$h \frac{u}{a} + k \frac{v}{b} + l \frac{w}{c} = 0$$

and $d(h, k, l)$ is the length of the perpendicular from the origin on to the plane

$$h \frac{u}{a} + k \frac{v}{b} + l \frac{w}{c} = 1.$$

It follows from 11.3.1, 11.3.2 and 11.3.3 that:

$$p_3(r) = \frac{1}{V} \sum_{h,k,l} \sum \sum |F(h, k, l)|^2 \frac{\sin 2\pi r/d(h, k, l)}{2\pi r/d(h, k, l)} \dots (11.3.4)$$

In a similar manner it is sometimes necessary to obtain radial distribution functions for the two dimensional series:

$$P(u, v) = \frac{1}{A} \sum_{h,k} |F(h, k, 0)|^2 \cos 2\pi \left(h \frac{u}{a} + k \frac{v}{b} \right) \quad \dots (11.3.5)$$

and it is easily shown that the relevant function is defined by:

$$2\pi r p_2(r) dr = \int_S P(u, v) dS \quad \dots (11.3.6)$$

where S is now a circle in the (u, v) plane with centre $(0, 0)$. The integration is performed as before to give:

$$p_2(r) = \frac{1}{A} \sum_{h,k} |F(h, k, 0)|^2 J_0^2 2\pi r / d(h, k, 0) \quad \dots (11.3.7)$$

where $J_0(z)$ is the zeroth order Bessel function.

Other types of syntheses can be obtained by projecting the density $\rho(x, y, z)$ contained between selected values of z , say, on to a plane parallel to (x, y) . These, whilst of considerable value to the specialist crystallographer, are not of general interest and require no novelties of technique, the interested reader is referred to the literature.⁽¹⁰⁾

11.4 FOURIER ANALYSIS

The converse process to Fourier synthesis may be expressed, in the one dimensional case, by

$$\rho(x) = \frac{1}{2} A_0 + \sum_{h=1}^{\infty} A_h \cos 2\pi \left(h \frac{x}{a} \right) + B_h \sin 2\pi \left(h \frac{x}{a} \right) \quad \dots (11.4.1)$$

where $\rho(x)$ is known in the repetition interval $(0, a)$ and it is required to determine the values of the Fourier coefficients (A_h, B_h) .

Formally we have at once from the orthogonality of the sine and cosine functions:

$$A_h = \frac{2}{a} \int_0^a \rho(x) \cos 2\pi \left(h \frac{x}{a} \right) dx \quad (h = 0, 1 \dots \infty) \quad \dots (11.4.2)$$

$$B_h = \frac{2}{a} \int_0^a \rho(x) \sin 2\pi \left(h \frac{x}{a} \right) dx \quad (h = 1, 2 \dots \infty). \quad \dots (11.4.3)$$

To evaluate the integrals practically we make use of the Euler-Maclaurin integration formula (Chapter 4, section 4, equation 4.4.5).

When it is noted that $\rho(0) = \rho(a)$ and that by definition the same is true for all derivatives so that $\rho^{(n)}(0) = \rho^{(n)}(a)$ the Euler-Maclaurin formula becomes in this case:

$$\begin{aligned} \int_0^{n\delta x} \rho(x) \cos 2\pi \left(h \frac{x}{a} \right) dx = \\ \frac{\delta x}{2} \left[\rho(0) + 2\rho(\delta x) \cos 2\pi \left(h \frac{\delta x}{a} \right) + 2\rho(2\delta x) \cos 2\pi \left(h \cdot 2 \frac{\delta x}{a} \right) + \dots \right. \\ \left. + 2\rho\{(n-1)\delta x\} \cos 2\pi \left\{ h(n-1) \frac{\delta x}{a} \right\} + \rho(0) \right] \\ = \delta x \sum_{r=0}^{n-1} \rho(r\delta x) \cos 2\pi \left(hr \frac{\delta x}{a} \right) \quad \dots (11.4.4) \end{aligned}$$

since $n\delta x = a$.

Similarly:

$$\int_0^{n\delta x} \rho(x) \sin 2\pi \left(h \frac{x}{a} \right) dx = \delta x \sum_{r=0}^{n-1} \rho(r\delta x) \sin 2\pi \left(hr \frac{\delta x}{a} \right) \quad \dots (11.5.5)$$

and both of these formulae are *exact*.

Whence, from 11.4.2 and 11.4.3

$$\left. \begin{aligned} A_h &= \frac{2}{n} \sum_{r=0}^{n-1} \rho(r\delta x) \cos 2\pi \left(hr \frac{\delta x}{a} \right) \\ B_h &= \frac{2}{n} \sum_{r=0}^{n-1} \rho(r\delta x) \sin 2\pi \left(hr \frac{\delta x}{a} \right) \end{aligned} \right\} \quad \dots (11.5.6)$$

it should be noticed that by virtue of their derivation these formulae are significant only if $h < n/2$.

The practical computation of Fourier coefficients can be achieved by the use of the Beevers-Lipson strips in a manner similar to that used for Fourier synthesis. This follows at once if it is noticed that equations 11.5.6 are identical with 11.1.2 if the (x) of the latter is replaced by (h) and the (h) by $\left(r \frac{\delta x}{a} \right)$. A difficulty arises with the normal strips because they extend to frequencies (h values) of only 20, but this causes trouble only when Fourier coefficients of order greater than 10 are required. For higher values the equivalent of the strips are available on punched cards to frequencies of 60.

The process of analysis can be extended to multi-dimensional Fourier syntheses of the type in equation 11.1.1, the results are obtained in an analogous manner but are rather cumbersome and will not be tabulated here.⁽¹¹⁾

REFERENCES

- ⁽¹⁾ BRAGG, W. H. and BRAGG, W. L., 'The Crystalline State,' vol. 1, p. 15, Bell (1933)
- ⁽²⁾ BEEVERS, C. A. and LIPSON, H., *Phil. Mag.*, 17 (1934), 855
- ⁽³⁾ — — *Proc. Phys. Soc.*, 48 (1936) 772
- ⁽⁴⁾ BOOTH, A. D., 'Fourier Technique in X-ray Organic Structure Analysis,' p. 55, Cambridge (1948)
- ⁽⁵⁾ LIPSON, H. and COCHRAN, W., 'The Crystalline State,' vol. 3, Bell (1954)
- ⁽⁶⁾ BOOTH, A. D., *Trans. Faraday Soc.*, 42 (1946), 444
- ⁽⁷⁾ — *ibid.*, 42 (1946), 617
- ⁽⁸⁾ — 'Fourier Technique in X-ray Organic Structure Analysis,' p. 18, Cambridge (1948)
- ⁽⁹⁾ WIENER, N., 'The Fourier Integral,' Cambridge (1932)
- ⁽¹⁰⁾ BOOTH, A. D., *Trans. Faraday Soc.*, 41, (1945), 434
- ⁽¹¹⁾ MACGILLAVRY, C. H. and PEPINSKY, R., 'Computing Methods and the Phase Problem,' p. 310, Penn. State College (1952)

INTEGRAL EQUATIONS

12.1 CLASSIFICATION

THE integral equations which have been studied most extensively by analysts⁽¹⁾ are those associated with the names of Fredholm and of Volterra, they differ only in the fact that the limits of Volterra's equation contain the independent variable.

Fredholm's equation is usually written in the two forms:

$$\int_a^b k(x, y)f(y)dy = g(x) \quad \dots (12.1.1)$$

$$\int_a^b k(x, y)f(y)dy = g(x) + f(x) \quad \dots (12.1.2)$$

and the method of solution is closely governed by the presence or absence of the wanted function, $f(x)$, on the right-hand side of the equation. The function $k(x, y)$ is known as the 'kernel' and, just as is the case in the classical analytical theory, the technique of solution is simplified when $k(x, y)$ is symmetric in x and y . Associated with the Fredholm equations is the eigenvalue problem:

$$\lambda \int_a^b k(x, y)f(y)dy = f(x) \quad \dots (12.1.3)$$

where, not only is $f(x)$ unknown, but a solution exists only for discrete values of λ which have also to be determined.

The Volterra equations are

$$\int_a^x k(x, y)f(y)dy = g(x) \quad \dots (12.1.4)$$

$$\int_a^x k(x, y)f(y)dy = g(x) + f(x) \quad \dots (12.1.5)$$

and it is intuitively evident that their solution will resemble that of one-point boundary differential equations, while that of the Fredholm equations will be more akin to the two-point type.

An integral equation is said to be non-singular when both k and f , as well as the limits of integration, are finite and continuous.

Other integral equations derive from Fourier's theorem and from the Laplace transform; since they are more properly a part of analysis they will not be considered further here.

Two classical equations are those of Abel and of Schlömilch. Abel's equation is:

$$g(x) = \int_a^x \frac{f(y)}{(x-y)^\mu} dy \quad \begin{matrix} 0 < \mu < 1 \\ a \leq x \end{matrix} \quad \dots (12.1.6)$$

which can be shown⁽²⁾ to have the solution:

$$f(y) = \frac{\sin(\mu\pi)}{\pi} \frac{d}{dy} \int_a^y \frac{g(x) dx}{(y-x)^{1-\mu}} \quad \dots (12.1.7)$$

so that, when $g(x)$ is known, $f(y)$ can be found either analytically or by numerical integration.

Schlömilch's equation is:

$$g(x) = \frac{2}{\pi} \int_0^{\pi/2} f(x \sin \theta) d\theta \quad -\pi \leq x \leq \pi \quad \dots (12.1.8)$$

and the solution is⁽³⁾:

$$f(x) = g(0) + x \int_0^{\pi/2} g'(x \sin \theta) d\theta \quad \dots (12.1.9)$$

which is also amenable to treatment.

Equations occur in which both integrals and derivatives are present, and these are usually called integro-differential equations. The method of numerical solution depends upon the type of boundary conditions which have to be satisfied, but otherwise follows closely upon the lines which will be indicated in sections 12.2 and 12.3.

12.2 VOLTERRA'S EQUATION

In the first place we may note that an equation of the type 12.1.4 can always be reduced to one of type 12.1.5 by a single differentiation with respect to x if $k(x, x) \neq 0$, or by a series of such differentiations until a derivative is found such that $k^{(n)}(x, x) \neq 0$.

The solution now proceeds in two parts, first initial values are obtained by means of a Taylor series expansion:

$$f(a+x) = f(a) + xf'(a) + \frac{1}{2}x^2f''(a) + \dots$$

where, from equation 12.1.5, we have:

$$\begin{aligned} f(a) &= -g(a) \\ f'(a) &= -g'(a) + k(a, a)f(a) \\ f''(a) &= -g''(a) + k(a, a)f'(a) + k'(a, a)f(a) + k(a, a)f'(a) \\ &\dots \\ f^{(n)}(a) &= -g^{(n)}(a) + \sum_{r=0}^{n-1} \left[\frac{d^{n-1-r}}{dx^{n-1-r}} \{k^{(r)}(x, x)f(x)\} \right]_{x=a} \end{aligned} \quad \dots (12.2.1)$$

Next we take Gregory's formula for the integral which may be derived from the Euler-Maclaurin formula 4.4.5 by substituting for the derivatives in terms of forward and backward differences.

The Gregory formula is:

$$\begin{aligned} \int_a^{a+n\delta x} f(x) dx &= \\ &\delta x \left[\frac{1}{2}f(a) + f(a+\delta x) + \dots + f(a+(n-1)\delta x) + \frac{1}{2}f(a+n\delta x) \right] \\ &+ \delta x \left[\frac{1}{12}\Delta - \frac{1}{24}\Delta^2 + \frac{1}{720}\Delta^3 - \frac{1}{1680}\Delta^4 + \dots \right] f(a) \\ &- \delta x \left[\frac{1}{12}\nabla + \frac{1}{24}\nabla^2 + \frac{1}{720}\nabla^3 + \frac{1}{1680}\nabla^4 + \dots \right] f(a+n\delta x) \end{aligned} \quad \dots (12.2.2)$$

and it is seen that, by using it to a prescribed number of terms, the given integral can be expressed in the form:

$$\int_a^{a+n\delta x} f(x) dx = \sum_{r=0}^n \delta x A_r f(a+r\delta x) \quad \dots (12.2.3)$$

where the coefficients A_r depend only upon the number of differences to which 12.2.2 is taken.

By means of 12.2.3 the integral equation 12.1.5 may be written:

$$\delta x \sum_{r=0}^n A_r k(a+n\delta x, a+r\delta x) f(a+r\delta x) = g(a+n\delta x) + f(a+n\delta x)$$

or:

$$\begin{aligned} [1 - \delta x A_n k(a+n\delta x, a+n\delta x)] f(a+n\delta x) &= -g(a+n\delta x) \\ &+ \delta x \sum_{r=0}^{n-1} A_r k(a+n\delta x, a+r\delta x) f(a+r\delta x) \end{aligned} \quad \dots (12.2.4)$$

whence, if the first (n) values of $f(a+r\delta x)$ are known, a further value may be obtained. The process can now be repeated, using the values $f(a+n\delta x) \dots f(a+\delta x)$, and the solution thus continued.

Care must be taken that differences of sufficiently high order be included in equation 12.2.2 for the rate of growth of error to be adequately small in the range over which a solution is required. For an indication of how this is done, and also for suggestions as to methods of estimating the error in particular cases, the reader is referred to the excellent paper of Fox and GOODWIN⁽⁴⁾.

12.3 FREDHOLM'S EQUATION

The solution of the Fredholm equation of the first kind, 12.1.1, is complicated by the fact that solutions may not be possible when k and g are connected in certain ways. An obvious example of this arises when $k(x, y)$ has the form $\sum_{r=1}^n X_r(x) \cdot Y_r(y)$, in which case equation 12.1.1 becomes:

$$\sum_{r=1}^n C_r X_r(x) = g(x) \quad \dots (12.3.1)$$

where:

$$C_r = \int_a^b Y_r(y) f(y) dy. \quad \dots (12.3.2)$$

Here no solution is possible unless $g(x)$ has the form 12.3.1 but, in this case, any solution is possible for which equation 12.3.2 is satisfied. Fox and GOODWIN⁽⁶⁾ conclude that no really satisfactory numerical method exists for the accurate solution of equations of this kind.

For Fredholm's equation of the second kind, 12.1.2, the position is more satisfactory. The approach is again to represent the integral by a finite difference approximation of the type given in equations 12.2.2 and 12.2.3, but from here there exists a choice of two methods.

In the first a sufficient number of differences are retained in equation 12.2.2 to ensure that the reduced form 12.2.3 gives an adequate representation of the integral over the range considered, in the second method equation 12.2.2 is written in the form:

$$\int_a^{a+n\delta x} f(x) dx = \delta x [\frac{1}{2} f(a) + f(a + \delta x) + \dots + f\{a + (n-1)\delta x\} + \frac{1}{2} f(a + n\delta x)] + \mathcal{E} \quad \dots (12.3.3)$$

where \mathcal{E} is a correction term compounded of the differences given

in equation 12.2.2. The integral equation is now written in one of the forms:

$$\delta y \sum_{r=0}^n A_r k(x, a + r\delta y) f(a + r\delta y) = g(x) + f(x) \quad \dots (12.3.4)$$

or:

$$\begin{aligned} \delta y [\frac{1}{2} k(x, a) f(a) + k(x, a + \delta y) f(a + \delta y) + \dots \\ + k\{x, a + (n-1)\delta y\} f\{a + (n-1)\delta y\} \\ + \frac{1}{2} k(x, a + n\delta y) f(a + n\delta y)] + \mathcal{E}_x = g(x) + f(x). \end{aligned} \quad \dots (12.3.5)$$

We notice that $a + n\delta y = b$ and then proceed to solve the $(n+1)$ simultaneous linear equations, which result from taking $x = a, a + \delta x, \dots, a + n\delta x$ in 12.3.4, for the $(n+1)$ values of $f(x)$. This method has the disadvantage that it is necessary to keep sufficient differences in deriving equation 12.2.3 to ensure the required final accuracy and since the behaviour of the solution is in general unknown, this may lead to wasted effort.

Fox and GOODWIN⁽⁶⁾ have suggested a modification of their method for the solution of ordinary differential equations (see page 68). In the present instance they solve the equations which result from substituting $x = a, a + \delta x, \dots, a + n\delta x$ in 12.3.5, initially neglecting the corrections \mathcal{E}_x . The set of approximations $f^{(0)}(a), f^{(0)}(a + \delta x), \dots, f^{(0)}(a + n\delta x)$ are then used to estimate the values of \mathcal{E}_x and the solution is re-computed to give a better approximation $f^{(1)}(a), f^{(1)}(a + \delta x)$, etc. This process is repeated until no significant change is produced between two successive approximations.

Alternative methods which have been suggested for the solution of equations of the Fredholm type include the deferred approach to the limit technique of RICHARDSON⁽⁷⁾, and the analogue of Picard's method for ordinary differential equations (see page 57) which sets up an iterative sequence:

$$f^{(n)}(x) = -g(x) + \int_a^b k(x, y) f^{(n-1)}(y) dy. \quad \dots (12.3.6)$$

Since it is necessary to calculate the value of each iterate at a sufficient number of points to enable the next integrations to be carried out to an adequate accuracy the method is likely to be tedious.

12.4 THE EIGENVALUE PROBLEM

The solution of the eigenvalue problem typified by equation 12.1.3 follows closely upon the lines just suggested for solving the Fredholm

equation 12.1.2. The equation is put into one of the finite difference forms:

$$\sum_{r=0}^n A_r k(x, a + r\delta y) f(a + r\delta y) = f(x)/\lambda \delta y \quad \dots (12.4.1)$$

or:

$$\begin{aligned} & \frac{1}{2} k(x, a) f(a) + k(x, a + \delta y) f(a + \delta y) + \dots \\ & + k(x, a + (n-1)\delta y) f(a + (n-1)\delta y) \\ & + \frac{1}{2} k(x, a + n\delta y) f(a + n\delta y) = f(x)/\lambda \delta y - \mathcal{E}/\delta y \\ & \dots (12.4.2) \end{aligned}$$

and the sets of simultaneous equations which result from taking $x = a, a + \delta x, a + 2\delta x, \dots, a + n\delta x$ are then used as a basis of computation.

For the equations which derive from 12.4.1 the values of λ and their associated vectors are determined in the manner described in Chapter 7 section 7.6 since it is evident that non-zero solutions exist only if the determinant:

$$\begin{vmatrix} A_0 k(0, 0) - 1/\lambda \delta y & A_1 k(0, 1) & \dots & A_n k(0, n) \\ A_0 k(1, 0) & A_1 k(1, 1) - 1/\lambda \delta y & \dots & A_n k(1, n) \\ \dots & \dots & \dots & \dots \\ A_0 k(n, 0) & A_1 k(n, 1) & \dots & A_n k(n, n) - 1/\lambda \delta y \end{vmatrix}$$

is zero. We have used the shorthand notation:

$$k(r, s) = k(a + r\delta x, a + s\delta y).$$

This technique pre-supposes a knowledge of the behaviour of the solution to enable appropriate values of the coefficients A_r to be chosen. The second method first obtains rough values of λ under the assumption that the corrections \mathcal{E}_x are zero and then uses each of these values of λ to derive a trial vector $f(a), f(a + \delta x), \dots, f(a + n\delta x)$. From this trial vector the corrections \mathcal{E}_x can be computed and then, in turn, an improved value of λ derived; the latter operation is best carried out by means of the second order process given in Chapter 7 section 6 equation 7.6.2.

The whole process is a cyclic one which is terminated as soon as the change between two steps is less than the desired value.

12.5 MONTE CARLO METHODS

Since the methods of solution for integral equations which we have just outlined aim at reducing the equation to a set of linear simultaneous equations, it is natural to suppose that, just as is the case for the latter, there will be Monte Carlo methods for solving integral equations.

The Fredholm equation arises in the study of the behaviour of neutrons which are incident on a plate of material in which collisions can give rise either to more or to less than one particle. This suggests⁽⁸⁾ the following Monte Carlo process for solving an equation of the type given in 12.1.2; first the equation is normalized so that:

$$\int_a^b g(x) dx = 1 \quad \text{and} \quad \int_{-\infty}^{\infty} k(x, y) dy \leq 1 \quad \dots (12.5.1)$$

$f(x)$ now represents the density of collisions in a population of particles in (a, b) when it is assumed that the density of first collisions is $g(x)$. In normalizing $k(x, y)$ a numerical factor, λ , will have been introduced and the equation can be written:

$$f(x) = g_n(x) + \lambda \int_a^b k_n(x, y) f(y) dy \quad \dots (12.5.2)$$

where the subscript (n) indicates that the function has been normalized in accord with 12.5.1, and the number λ is interpreted to be the mean number of particles which remain after each collision.

To solve the problem, we first pick a random number r_0 from a population which is uniformly distributed in $(0, 1)$, a first collision position, x_1 , is then determined from the equation:

$$\int_a^{x_1} g(x) dx = r_0 \quad \dots (12.5.3)$$

and a record is made that one collision has occurred at x_1 . Now λ will not, in general, be a whole number, so that to determine the number of particles which result from the collision, a statistical procedure must be used. A simple method is to take the number of particles to be that integer just less than or just greater than λ according to whether a second random number, r_1 , is less than, or greater than, the fractional part of λ .

Each particle which results from the first collision is now followed, and second collision positions, x_2 , are determined from:

$$\int_{-\infty}^{x_2} k(x_1, y) dy = r_2 \quad \dots (12.5.4)$$

where r_2 is again a random number in $(0, 1)$.

When $r_2 > \int_a^b k(x_1, y) dy$ the particle is assumed to have escaped from (a, b) and is no longer of interest; when all of the particles which resulted from the first collision have been followed to their termination (either by attenuation or by escape) a new particle is started and the process repeated. After a number of repetitions of this process a graphical record will be available on which is marked each collision which occurred in (a, b) ; a smooth curve drawn through this gives the required approximation to $f(x)$.

In a practical application of the above process a sample of 100 starting neutrons yielded a final accuracy of about 10 per cent. It is perhaps worth mentioning that equations 12.5.3 and 12.5.4 are normally solved by means of previously constructed graphs.

REFERENCES

- (1) WHITTAKER, E. T. and WATSON, G. N., 'Modern Analysis,' p. 211, 4th Edn. Cambridge (1927)
- (2) — — *ibid.*, p. 229
- (3) — — *ibid.*, p. 229
- (4) FOX, L. and GOODWIN, E. T., *Phil. Trans. Roy. Soc., A*, 245 (1953) 524
- (5) — — *ibid.*, p. 517
- (6) — — *ibid.*, p. 503
- (7) RICHARDSON, L. F. and GAUNT, J. A., *Phil. Trans. Roy. Soc., A*, 226 (1927) 299
- (8) SPINRAD, B. I., GOERTZEL, G. H. and SNYDER, W. S., 'Monte Carlo Method,' p. 4 *Nat. Bur. Stand. Appl. Math. Series No. 12*, Washington (1951)

SELECT BIBLIOGRAPHY

This is not intended to be exhaustive, but merely to indicate recent key publications containing references which will enable the reader to trace the previous literature of any particular branch of the subject.

General

- BEREZIN, I. S. and ZHIDKOV, N. P., 'Computing Methods,' (2 vols.), Pergamon, Oxford, 1965
- BUTLER, R. and KERR, E., 'An Introduction to Numerical Mathematics,' Pitman, London, 1962
- FREEMAN, H., 'An elementary treatise on Actuarial Mathematics,' Ch. II—VIII, XVII, Cambridge University Press, London, 1931
- FROBERG, C. E., 'Introduction to Numerical Analysis,' Addison-Wesley, Cambridge, Mass., 1965
- HAMMING, R. W., 'Numerical Mathematics for Scientists and Engineers,' McGraw-Hill, New York, 1962
- HARTREE, D. R., 'Numerical Analysis,' 2nd edn., Oxford University Press, London, 1957
- HENRICI, P., 'Elements of Numerical Analysis,' Wiley, New York, 1964
- HILDEBRAND, F. B., 'Introduction to Numerical Analysis,' McGraw-Hill, New York, 1956
- HOUSEHOLDER, A. S., 'Principles of Numerical Analysis,' McGraw-Hill, New York, 1953
- KOPAL, Z., 'Numerical Analysis,' Chapman & Hall, London, 1955
- KUO, S. S., 'Numerical Methods and Computers,' Addison-Wesley, Mass., 1965
- LANCE, G. N., 'Numerical Methods for High-Speed Computers,' Iliffe, London, 1960
- MILNE, W. E., 'Numerical Calculus,' Princeton University Press, 1949
- NIELSEN, K. L., 'Methods in Numerical Analysis,' Macmillan, New York, 1956
- NORKIN, S. B., 'The Elements of Computational Mathematics,' Pergamon, Oxford, 1965
- RALSTON, A. and WILF, H. S. (eds.), 'Mathematical Methods for Digital Computers,' Wiley, New York, 1960
- REDISH, K. A., 'An Introduction to Computational Methods,' English Universities Press, London, 1961
- SALVADORI, M. G. and BARON, M. L., 'Numerical Methods in Engineering,' Prentice Hall, New Jersey, 1964
- SCARBOROUGH, J. B., 'Numerical Mathematical Analysis,' 5th edn., Johns Hopkins Press, Baltimore, 1962
- STANTON, R. G., 'Numerical Methods for Scientists and Engineers,' Prentice Hall, New Jersey, 1961
- 'Modern Computing Methods,' 2nd edn., H.M. Stationery Office, London, 1961
- WHITTAKER, E. T. and ROBINSON, G., 'The Calculus of Observations,' 4th edn., Blackie, London, 1949

Finite Differences

- BOOLE, G., 'A Treatise on the Calculus of Finite Differences,' 2nd edn., Macmillan, London, 1872
- DAVIS, P. J., 'Interpolation and Approximation,' Blaisdell, New York, 1963
- FORT, T., 'Finite Differences,' Oxford University Press, London, 1948
- JORDAN, C., 'Calculus of Finite Differences,' Chelsea, New York, 1947
- MILNE-THOMSON, L. M., 'The Calculus of Finite Differences,' Macmillan, London, 1951
- STEFFENSON, J. F., 'Interpolation,' Williams & Wilkins, Baltimore, 1927

SELECT BIBLIOGRAPHY

Integration and the Solution of Differential Equations

- ALLEN, D. N. de G., 'Relaxation Methods,' McGraw-Hill, New York, 1954
 BENNETT, A. A., MILNE, W. E. and BATEMAN, H., 'Numerical Integration of Differential Equations,' *Bull. Nat. Res. Coun. Wash.*, 92 (1933) 51
 COLLATZ, L., 'Numerische Behandlung von Differentialgleichungen,' Springer, Berlin, 1955
 FORSYTHE, G. E. and WASOW, W. R., 'Finite Difference Methods for Partial Differential Equations,' Wiley, New York, 1960
 FOX, L., 'The Numerical Solution of Two-Point Boundary Problems in Ordinary Differential Equations,' Oxford University Press, London 1957
 — 'Numerical Solution of Ordinary and Partial Differential Equations,' Addison-Wesley, Cambridge, Mass., 1962
 GIBB, D., 'Interpolation and Numerical Integration,' London, 1915
 LEVY, H. and BAGGOTT, E. A., 'Numerical Studies in Differential Equations,' vol. 1, Watts, London, 1934
 LOWAN, A. N., 'Tables of Functions and Zeros of Functions,' *Nat. Bur. Stand. Appl. Maths. Series*, No. 37, Washington, 1954
 MILNE, W. E., 'Numerical Solution of Differential Equations,' Wiley, New York, 1953
 SHAW, F. S., 'An Introduction to Relaxation Methods,' Dover, New York, 1953
 SOUTHWELL, R. V., 'Relaxation Methods in Engineering Science,' Oxford University Press, London, 1940
 — 'Relaxation Methods in Theoretical Physics,' Oxford University Press, London, 1946
 TITCHMARSH, E. C., 'Eigenvalue Expansions,' Vol. 1, Oxford University Press, London, 1946; Vol. 2, Oxford University Press, London, 1958

Solution of Simultaneous Linear Equations and Inversion of Matrices

- FADDEEV, D. K. and FADDEEVA, V. N., 'Computational Methods in Linear Algebra,' Freeman, London, 1963
 FOX, L., 'Introduction to Numerical Linear Algebra,' Oxford University Press, London, 1964
 HOUSEHOLDER, A. S., 'The Theory of Matrices in Numerical Analysis,' Blaisdell, New York, 1964
 PAIGE, L. J. and TAUSSKY, O. (Ed.), 'Simultaneous Linear Equations and the Determination of Eigenvalues,' *Nat. Bur. Stand. Appl. Maths. Series*, No. 29, Washington, 1953. An invaluable bibliography of over 450 titles, together with several original papers
 TAUSSKY, O. (Ed.), 'Contributions to the Solution of Systems of Linear Equations and the Determination of Eigenvalues,' *Nat. Bur. Stand. Appl. Maths. Series*, No. 39, Washington, 1954
 WILKINSON, J. H., 'The Evaluation of Latent Roots and Vectors of Matrices on the Pilot Model of the A.C.E.,' *Proc. Camb. Phil. Soc.*, 50 (1954) 536
 — 'Rounding Errors in Algebraic Processes,' Prentice-Hall, New Jersey, 1963
 — 'The Algebraic Eigenvalue Problem,' Oxford University Press, London, 1965

Non-linear Algebraic Equations

- BOOTH, A. D., 'An Application of the Method of Steepest Descents,' *Quart. J. Mech.*, 2 (1949) 460
 OLVER, F. W. J., 'The Evaluation of Zeros of High-degree Polynomials,' *Phil. Trans.*, A, 244 (1952) 385
 TRAUB, J. F., 'Iterative Methods for the Solution of Equations,' Prentice-Hall, New Jersey, 1964

SELECT BIBLIOGRAPHY

Fourier Analysis and Synthesis

- BARBER, N. F., 'Experimental Correlograms and Fourier Transforms,' Pergamon, Oxford, 1961
 BOOTH, A. D., 'Fourier Technique,' Cambridge University Press, London, 1948
 MANLEY, R. G., 'Waveform Analysis,' Chapman & Hall, London, 1945

Polynomial Expansions

- CLENSHAW, C. W., 'Chebyshev Series for Mathematical Functions,' H.M. Stationery Office, London, 1965
 HASTINGS, C., 'Approximations for Digital Computers,' Princeton University Press, 1955
 LANCZOS, C., 'Tables of Chebyshev Polynomials,' *Nat. Bur. Stand. Appl. Maths. Series* No. 9, Washington, 1952
 POUSSIN, C. DE LA VALLÉE, 'Sur la méthode de l'approximation minimum,' *Soc. Sci. de Bruxelles, Annales, Second Partie, Mémoires*, 35 (1911) 1
 SZEGÖ, G., 'Orthogonal Polynomials,' *Amer. Math. Soc. Colloq. Publ.* vol. XXIII, New York, 1939

Integral Equations

- COURANT, R. and HILBERT, D., 'Methods of Mathematical Physics,' vol. 1, Interscience, New York, 1953
 FOX, L. and GOODWIN, E. T., 'The Numerical Solution of Non-singular Integral Equations,' *Phil. Trans.*, A, 245 (1953) 501
 LONSETH, A. T., 'Approximate Solutions of Fredholm-type Integral Equations,' *Bull. Amer. Math. Soc.*, 60 (1954) 415
 MUSKHELISHVILI, N. I., 'Singular Integral Equations,' Noordhoff, Groningen, 1953
 WHITTAKER, E. T. and WATSON, G. N., 'Modern Analysis,' 4th edn., Cambridge University Press, 1927
 YOUNG, A., 'Approximate Product-integration,' *Proc. Roy. Soc. A*, 224 (1954) 552-561
 — 'The Application of Product-integration to the Numerical Solution of Integral Equations,' *Proc. Roy. Soc. A*, 224 (1954) 561-573

Monte Carlo Methods

- CASHWELL, E. D. and EVERETT, C. J., 'Monte Carlo Methods for Random Walk Problems,' Pergamon, Oxford, 1959
 CURTISS, J. H., 'Monte Carlo Methods for the Iteration of Linear Operators,' *Nat. Bur. Stand. Rep. No.* 2365, Washington, 1953
 HAMMERSLEY, J. M. and HANDSCOMB, D. C., 'Monte-Carlo Methods,' Methuen, London, 1964
 — and MORTON, K. W., 'A New Monte Carlo Technique: Antithetic Variates,' *Proc. Camb. Phil. Soc.*, 52 (1956) 449
 HOUSEHOLDER, A. S. (Ed.), 'Monte Carlo Method,' *Nat. Bur. Stand. Appl. Maths. Series* No. 12, Washington, 1951
 TOCHER, K. D., 'Introduction to Monte Carlo Methods,' English Universities Press, London, 1963

Automatic Digital Calculators and Programming

- BERKELEY, E. C. and WAINWRIGHT, L., 'Computers, Their Operation and Applications,' Reinhold, New York, 1956
 BOOTH, A. D., 'Computers in Action,' Pergamon, Oxford, 1965
 — 'Computers and Automation,' Staples Press, London, 1965
 — and BOOTH, K. H. V., 'Automatic Digital Calculators,' 3rd edn., Butterworths, London, 1965

SELECT BIBLIOGRAPHY

- BOOTH, K. H. V., 'Programming for an Automatic Digital Calculator,' Butterworths, London, 1958
 COLEBROOK, F. M. (Ed.), 'Automatic Digital Computation,' H. M. Stationery Office, London, 1954
 HARTREE, D. R., 'Calculating Instruments and Machines,' Cambridge, 1950
 STIFLER, W. W. (Ed.), 'High Speed Computing Devices,' McGraw-Hill, New York, 1950
 WARE, W., 'Digital Computer Technology and Design,' (2 vols.), Wiley, New York, 1963

Mathematical Tables

- COMRIE, L. J., 'Chambers Four-figure Mathematical Tables,' Chambers, Edinburgh, 1947
 ERDÉLYI, A., MAGNUS, W., OBERHETTINGER, F. and TRICOMI, F. G., 'Higher Transcendental Functions,' 3 vols., McGraw-Hill, New York, 1953
 FLETCHER, A., MILLER, J. C. P. and ROSENHEAD, L., 'Index of Mathematical Tables,' Scientific Computing Service, London, 1946
 FLÜGGE, W., 'Four-place Tables of Transcendental Functions,' Pergamon, London, 1954
 JAHNKE, E. and EMDE, F., 'Tables of Higher Functions,' Teubner, Leipzig, 1952

NAME INDEX

- ABEL, N. H., 198
 ADAMS, J. C., 2, 59, 64, 65, 67, 70
 AL-KHOWARIZIMI, 1
 ALLEN, D. N. DE G., 206
 ARCHIMEDES, 1
 BABBAGE, C., 2
 BAGGOT, E. A., 74, 206
 BAIRSTOW, L., 168, 170
 BARGMAN, V., 123
 BARLOW, P., 3
 BASHFORTH, F., 59, 65, 67, 70
 BATEMAN, H., 206
 BEEVERS, C. A., 187, 188, 189, 190, 195, 196
 BENNETT, A. A., 206
 BERNOULLI, J., 36, 37
 BESSEL, F. W., 20, 22, 184, 185, 194
 BICKLEY, W. G., 32, 37, 53
 BLANCH, G., 127, 155
 BOOLE, G., 155, 205
 BOOTH, A. D., 6, 177, 196, 206, 207, 208
 BOOTH, K. H. V., 6, 208
 BRAGG, W. H., 196
 BRAGG, W. L., 187, 196
 BRIGGS, H., 2
 BROMWICH, T. J. I'A., 56, 58
 CARDAN, H., 162
 CAYLEY, A., 83
 CHAMBERS, W., 3
 CHEBYSHEV, P. L., 39, 42, 43, 49, 51, 183, 189, 185
 CHOLESKI, 81, 82
 CHRISTOFFEL, E. B., 43
 COCHRAN, W., 196
 COLEBROOK, F. M., 208
 COMRIE, L. J., 3, 19, 23, 24, 29
 COTES, R., 35, 36, 51, 67, 70
 COURANT, R., 53, 123, 125, 126, 136, 186, 207
 CURTISS, J. H., 155, 207
 DALE, J. B., 3
 DAVIDS, N., 46, 53
 DE COLMAR, C. X. T., 2
 DE HAHN, D. B., 3
 DE VOGELAERE, R., 71, 75
 DIOPHANTUS, 1
 DOOLITTLE, M. H., 7
 DWIGHT, H. B., 3
 EDELMAN, G. M., 155
 EMDE, F., 53
 ERDÉLYI, A., 208
 EULER, L., 36, 37, 55, 56, 57, 194, 195, 199
 EVERETT, J. D., 11, 21, 22, 63
 FERRAR, W. L., 123
 FERRARI, L., 162
 FILON, L. N. G., 50, 51, 53
 FLETCHER, A., 3, 208
 FLÜGGE, W., 208
 FORSYTH, G. E., 123
 FORT, 56, 205
 FOURIER, J., 127, 179, 180, 187, 194, 195, 196, 198
 FOX, L., 72, 75, 200, 201, 204, 206
 FREDHOLM, E. I., 197, 200, 201, 203
 FREEMAN, H., v, 20, 29, 205
 FRIEDRICHS, K., 125
 GAUNT, J. A., 204
 GAUSS, K. F., 19, 20, 43, 46, 50
 GIBB, D., 206
 GIVENS, W., 118, 123
 GOERTZEL, G. H., 204
 GOLDSTINE, H. H., 115, 123
 GOODWIN, E. T., 72, 75, 186, 200, 201, 204, 207
 GRÄEFFE, C. H., 162
 GRAM, J. P., 96, 110
 GREGORY, J., 17, 18, 34, 35, 199
 GUNTER, E., 2
 HAMILTON, W. R., 83
 HARTREE, D. R., 67, 74, 75, 127, 155, 177, 205, 208
 HERMITE, C., 48, 182, 183, 185
 HERON, 1
 HEUN, K., 64
 HICKS, B. L., 155
 HILBERT, D., 53, 123, 186, 207
 HOBSON, E. W., 53
 HORNER, W. G., 162
 HOUSEHOLDER, A. S., 118, 123, 206, 207
 JACOBI, K. G. J., 115, 118
 JAHNKE, E., 3, 53
 JENNINGS, J. C. E., 72, 75
 JORDAN, C., 205

NAME INDEX

KOPAL, Z., 25, 29, 53
 KRONECKER, L., 179
 KUTTA, W., 64, 65

 LAGRANGE, J., 27, 47
 LAGUERRE, E., 47, 48, 182, 183, 185
 LANCZOS, C., 53, 118, 186, 207
 LANDAU, H. G., 155
 LAPLACE, P. S., 124, 128, 129, 131, 139, 140, 145, 153, 198
 LEIBLER, R. A., 123
 LEVENSON, A., 46, 53
 LE VERRIER, U. J. J., 2
 LEVY, H., 74, 206
 LEWY, H., 125
 LIDSTONE, G. J., 29
 LIN, S. N., 168, 169
 LIPSON, H., 187, 188, 189, 190, 195, 196
 LOBATTO, 50
 LONSETH, A. T., 207
 LOTKIN, M., 155
 LOWAN, A. N., 46, 53

 MACGILLAVRY, C. H., 196
 MACLAURIN, C., 36, 37, 39, 51, 57, 61, 194, 195, 199
 MAGNUS, W., 208
 MANLEY, R. G., 207
 MERSENNE, P., 2
 MILLER, J. C. P., 3, 208
 MILNE, W. E., 51, 53, 59, 67, 70, 74, 75, 126, 131, 132, 155, 205, 206
 MILNE-THOMSON, L. M., 39, 53, 205
 MORRIS, J., 101, 104, 123
 MURRAY, F. J., 123
 MUSKHELISHVILI, N. I., 207

 NAPIER, J., 2
 NEWTON, I., 1, 17, 18, 34, 35, 36, 51, 52, 67, 70, 164, 168, 191

 OBERHETTINGER, F., 208
 OLVER, F. W. J., 177, 206

 PAIGE, L. J., 123, 206
 PEPINSKY, R., 196
 PIAGGIO, H. T. H., 64
 PICARD, C. E., 60, 61, 62, 63, 201
 PIERCE, B. O., 3
 POISSON, S. D., 114, 140, 153, 154
 PYTHAGORAS, 1

RADAU, R., 50
 RAPHSOON, J., 164, 166, 169, 191
 REIZ, A., 50
 RICHARDSON, L. F., 101, 123, 201, 204
 ROBINSON, G., v, 59, 74, 177, 205
 RODRIGUES, O., 45
 ROSENHEAD, L., 3, 208
 RUNGE, C., 64, 65

 SCARBOROUGH, J. B., 205
 SCHLÖMILCH, O., 198
 SCHMIDT, E., 96, 110
 SHAPIRO, A. H., 155
 SHAW, F. S., 206
 SIMPSON, T., 34, 39, 51, 67, 70
 SNYDER, W. S., 204
 SOMMERFELD, A., 155
 SOUTHWELL, R. V., 138, 147, 155, 206
 SPINRAD, B. I., 204
 STEFFENSON, J. F., 205
 STIFLER, W. W., 208
 STIRLING, J., 19, 20
 SZEGÖ, G., 207

 TARTAGLIA, N., 162
 TAUSSKY, O., 123, 206
 TAYLOR, B., 61, 68, 72, 164, 166, 171, 198
 TICHMARSH, E. C., 75, 206
 TRICOMI, F. G., 208

 ULAM, S., 121

 VLACQ, A., 2
 VOLTERRA, V., 197, 198
 VON NEUMANN, J., 115, 121, 123

 WASOW, W., 123
 WATSON, G. N., 204, 207
 WEDDLE, T., 35, 39
 WEFA, 1
 WHITTAKER, E. T., v, 59, 74, 177, 204, 205, 207
 WIENER, N., 196
 WILKINSON, J. H., 119, 121, 123, 172
 WOMERSLEY, J. F., 127, 155

 YOWELL, E. C., 132

SUBJECT INDEX

Advance to a finer net, in relaxation 142, 149
 Algebraic non-linear equations, 156
 Algol, 4
 Approximating functions, 178
 Arabia, 1
 Automatic digital calculator, 2, 4, 5, 42, 81, 113, 118, 132, 138, 145, 156, 178, 187
 Averaging operator, 14

 Babylon, 1
 Back substitution, 81
 Backward difference, 12
 Bairstow's method, 170
 Bashforth-Adams method, 65, 70
 Beavers-Lipson strips, 187, 195
 Bernoulli numbers, 36, 37
 Bessel, function, 184, 194
 interpolation formula, 20, 22, 24
 Block relaxation, 146
 Boundary conditions, 59, 136
 Brunsvig calculator, 2, 4

 Canonical form, 69, 70, 134, 135
 Casting out nines, 1
 Cayley-Hamilton theorem, 83
 Central difference, 11, 13, 21, 31,
 characteristic equation, 83, 101, 103
 numbers, 78, 90
 Characteristics, 134, 135
 Chebyshev integration formula, 39
 polynomials, 49, 183, 184
 Choleski method, 81
 Christoffel numbers, 43
 Closed integration formula, 35,
 Code, 4
 Conditions, of simultaneous equations, 84, 89, 108
 Contour maps, 160
 Convolution, 193
 Critical table, 10
 Curved boundaries, in relaxation, 144

 Detection of errors, 15
 Determinant, 77, 79
 Difference, backward, 12
 central, 11, 13, 21, 31
 divided, 25
 forward, 12
 leading, 7
 modified, 23

 Differential equation, solution of ordin-
 ary, 59
 partial, 124
 Fourier synthesis, 190
 Differentiation, numerical, 30, 31, 32
 Diffusion, equation of, 125
 Divided difference, 25
 Dyadic, 85, 91

 Egypt, 1
 Eigenvalue problem, for differential
 equations, 74
 for integral equations, 197, 201
 Eigenvalues, 74, 78, 90, 197, 201
 Electron lens, 72
 Elliptic integral, 72
 partial differential equations, 124,
 135, 136, 139, 145
 E.N.I.A.C., 2
 Errors, computational, 5
 detection of, 15
 experimental, 4
 in solution of Laplace equation, 153
 Escalator process, 101
 Euler-Maclaurin expansion, 36, 57, 194
 Euler transformation of series, 55, 56
 Everitt interpolation formula, 21, 22, 63

 Facit calculator, 4
 Factorial function, 17, 55
 False position, rule of, 1, 162
 Faltung, 193
 Finite difference, 7
 Forward difference, 12
 Fourier analysis, 194
 series, 127, 179, 189

 Gauss, backward interpolation formula,
 19
 forward interpolation formula, 19
 integration formula, 43, 46
 Generating function, 185
 Gram-Schmidt process, 96, 110
 Givens' method, 118
 Greece, 1
 Gregory integration formula, 199
 Group relaxation, 146

 Heat flow, equation of, 125
 Hermite polynomials, 48, 182, 183
 Hollerith punched card calculator, 2,
 132
 Homogeneous equations, 77

SUBJECT INDEX

Hyperbolic partial differential equation, 124, 127, 133, 135

India, 1

Integral equation, 197

Integrals, oscillating, 50

Integrals, tables of, 3

Integration, numerical, 30, 33, 38

Integro-differential equations, 198

Interpolation, 12, 17

Inverse interpolation, 27

matrix, 77

Iteration, method of, 164

Iterative process, 1, 162, 201

order of, 162, 192

Jacobi's method, 115

Jury problem, 59

Kernel, 197

Kronecker delta function, 123, 179

Lagrangean interpolation, 25, 44, 47

Laguerre polynomials, 47, 182, 183

Laplace equation, 124, 131, 139, 140, 153

operator, 128, 129, 153

Latent roots, 78, 90

vector, 78, 90

Leading difference, 7

Least square approximation, 179

Legendre polynomials, 45, 180, 181, 185

Linear interpolation, 9

Lin's method, 168

Lobatto's formula, 50

Marchant calculator, 4

Marching problem, 59

Matrix, 77, 137, 140

rotation, 115

Mercedes-Euclid calculator, 4

Mersenne number, 2

Mesh size, choice of, in relaxation, 141, 142, 146

Modified difference, 23

Monte Carlo method, 121, 153, 203

M.T.A.C., 3

Multiplicity correction, in multi-dimensional summation, 190

Neutron, 203

Newton-Cotes, integration formula, 35, 36, 67, 70

-Gregory, interpolation formula, 17, 18, 26, 34, 35

Newton-Raphson, iterative process, 164, 166, 169, 191

Non-linear algebraic equations, 156

Non-singular integral equation, 198

matrix, 77

Open integration formula, 36, 57, 52

Operators, 13, 30, 128, 153

Order, of an iterative process, 162, 192

Ordnance survey, 76

Orthogonal vectors, 91, 108

Orthogonality conditions, 109

Orthonormal functions, 179

transformation, 115

Oscillating integrals, 50

Parabolic partial differential equation, 124, 133, 135

Partial differential equation, 124

Pattern sensitivity, of a computer, 5

Patterson function, 193

Picard's method for initial values, 60, 201

Pivotal condensation, 81, 138

Poisson equation, 124, 140, 153, 154

Polynomial tabulation of, 5, 6

equation, solution of, 156

Positive definite form, 79, 105, 110, 111

Proportional parts, 8

Purification process, 101

Quadratic form, 79, 86, 101, 105, 124

Radau's formula, 50

Radial Fourier synthesis, 192, 193

Random walk, 122, 153, 154

Refinement, of approximate solutions, 156

Regula falsi, 1, 162

Relaxation, 59, 107, 138, 140, 145, 152, 171

Residual, 84, 152

Rodrigues' formula, 45

Rotation, of a matrix, 115

Runge-Kutta method, 64, 65

S.E.A.C., 2

Simpson's rule, 34, 39, 41, 53, 64, 67

Simultaneous linear equations, 76

Sine, calculation of, 5

Steepest descent, 105, 110, 114, 145, 175

Stirling's interpolation formula, 20

Summation of series, 54

Survey, to find approximate solution, 156

Tables, use of, 7

Tangent, 2

Tensor, 85, 91

Three eights rule, 35

SUBJECT INDEX

Throwback, 24

Trace, of a matrix, 84

Transition probability, 122

Transpose, of a matrix, 77

Under-determined sets of equations, 156

Vibrating string, equation of, 127

Wave equation, 132

Weddle's rule, 35, 39

Weight function, 43, 180, 182

Wilkinson's method, 119

X-ray crystallography, 160, 187, 192

NUMERICAL METHODS — BOOTH

THE
EDITION